



## ECE 538

### Advanced Computer Architecture

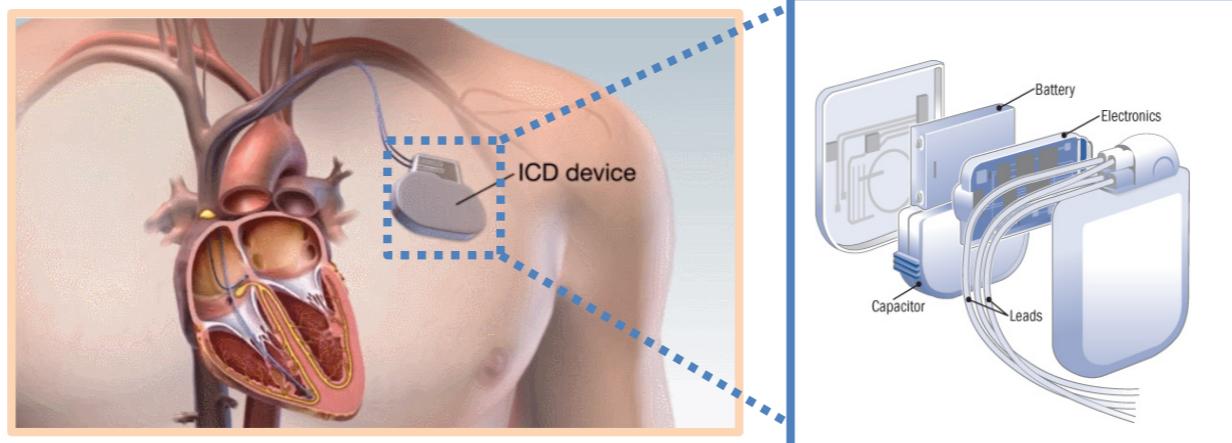
Instructor: Lei Yang

Department of Electrical and Computer Engineering

November 29, 2021

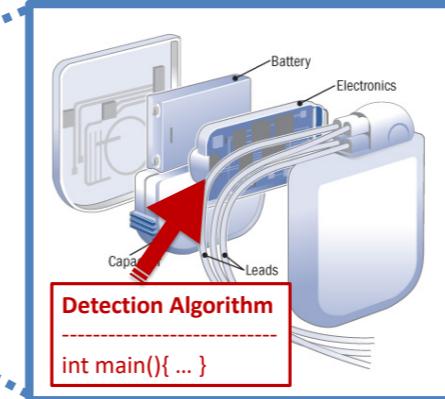
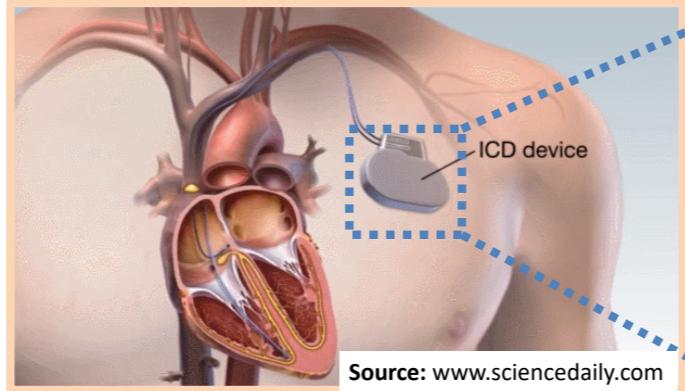
Hardware/Software Co-Exploration  
of Resource Constrained Edge AI Systems

- IMPLANTABLE CARDIOVERTER DEFIBRILLATORS (ICD)



Source: [www.sciencedaily.com](http://www.sciencedaily.com)

## ■ IMPLANTABLE CARDIOVERTER DEFIBRILLATORS (ICD)



### Problem in Existing Detection Approaches:



- Heuristic algorithms are not accurate
- More than **40% shocks** are inappropriate or missing

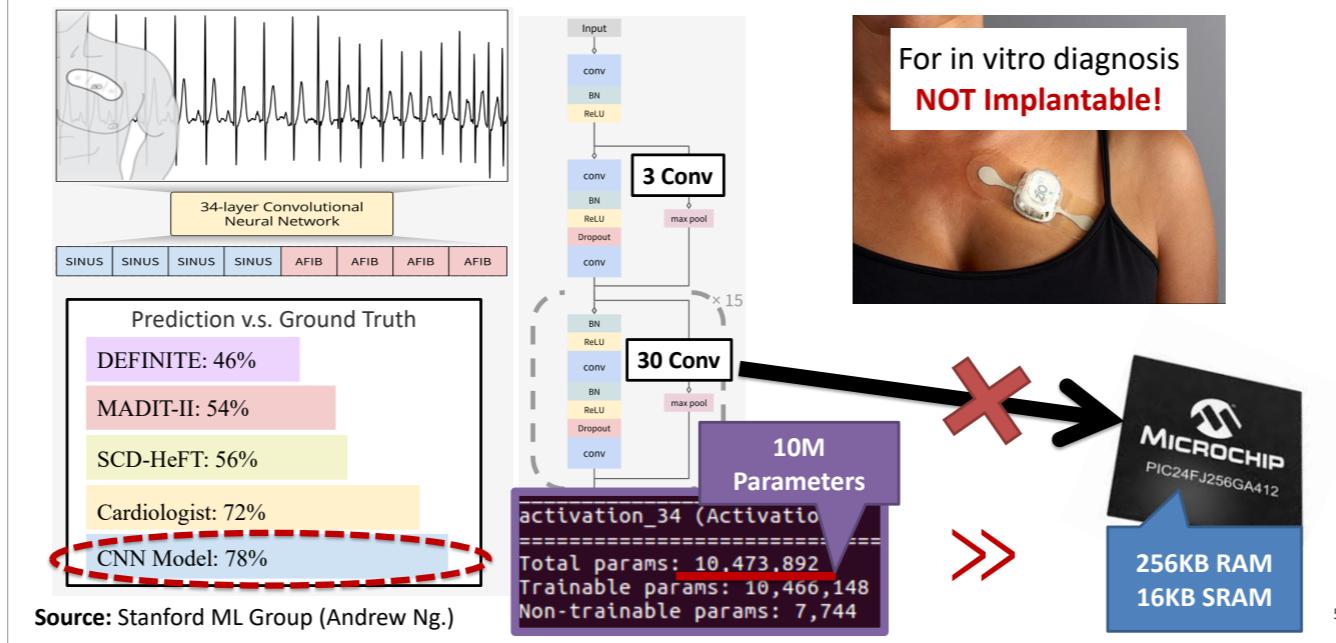
### Appropriate Shocks Ratio

DEFINITE: 46%

MADIT-II: 54%

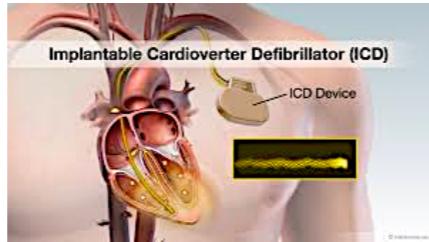
SCD-HeFT: 56%

- MACHINE LEARNING: BETTER THAN CARDIOLOGIST, BUT ...



5

## ■ PUT MACHINE LEARNING MODEL IN RESOURCE CONSTRAINED EDGE DEVICES



### VISION: ML in ICD Devices

- **Accuracy** --> Higher than 72% (Cardiologist)
- **Area** --> Less than  $3\mu\text{m}^2$
- **Power** --> Less than  $20\mu\text{W}$
- **Latency** --> Less than 10 heartbeat



### Voice language translation in AI Glasses



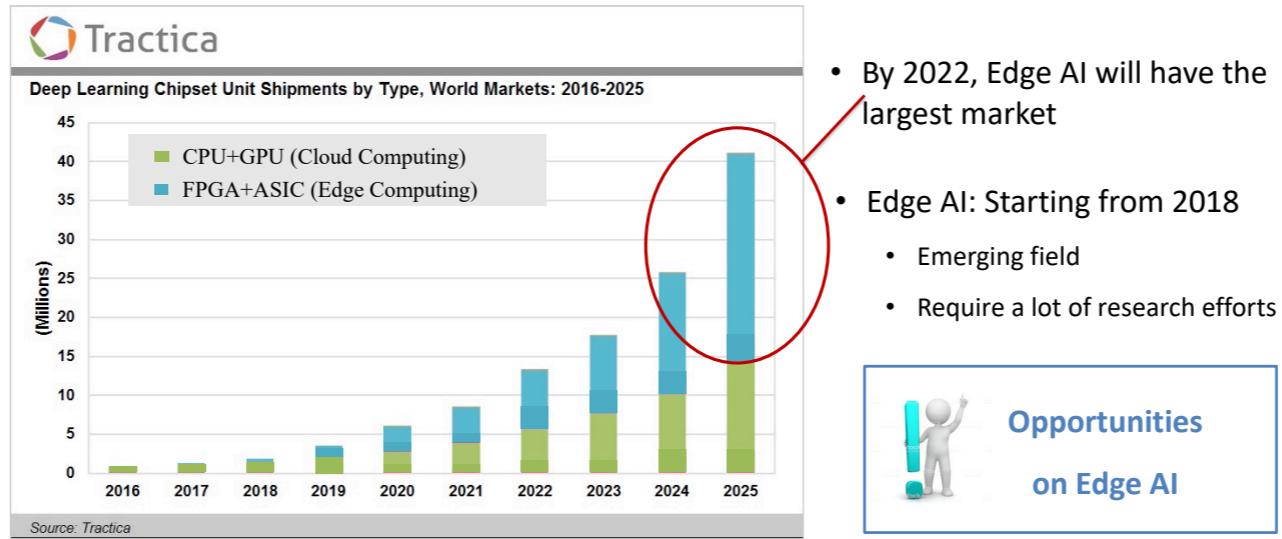
Surveillance Self-Driving Smart Home Crime Detection



6

Types or kinds of applications in H.AI especially during COVID-19.

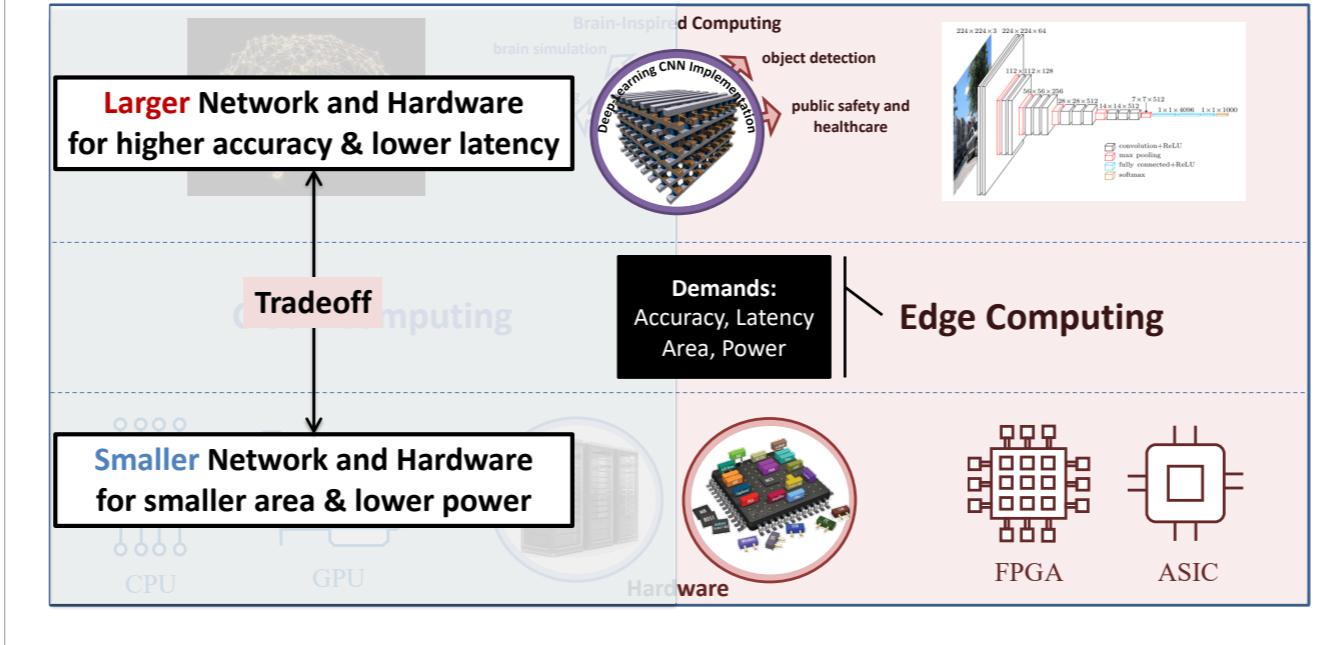
- EDGE AI WILL HAVE A LARGE MARKET



7

Number of applications increasing. By 2022 computing will have a larger market than cloud computing. Amount of data will make it expensive to move. Need to move to more real-time data. Privacy is a concern, but H.AI is still the future. Around April 2016 Intel et al. announced their Edge AI computing platform. Edge computing is the concept of doing all data work in close proximity. Real need of developers worldwide to train developers to deal with this scenario. Real need for AI to be completed on-device. App handles personal data instead of storing in cloud. Transferring the data over networks and then to cloud creates a security risk. SW engineers et al. have been working on AI applications, now have a new direction to take their work. More AI app development. These can be viewed from a more user-centric and is the next evolution of the AI revolution. However current state there are challenges including hardware and software design.

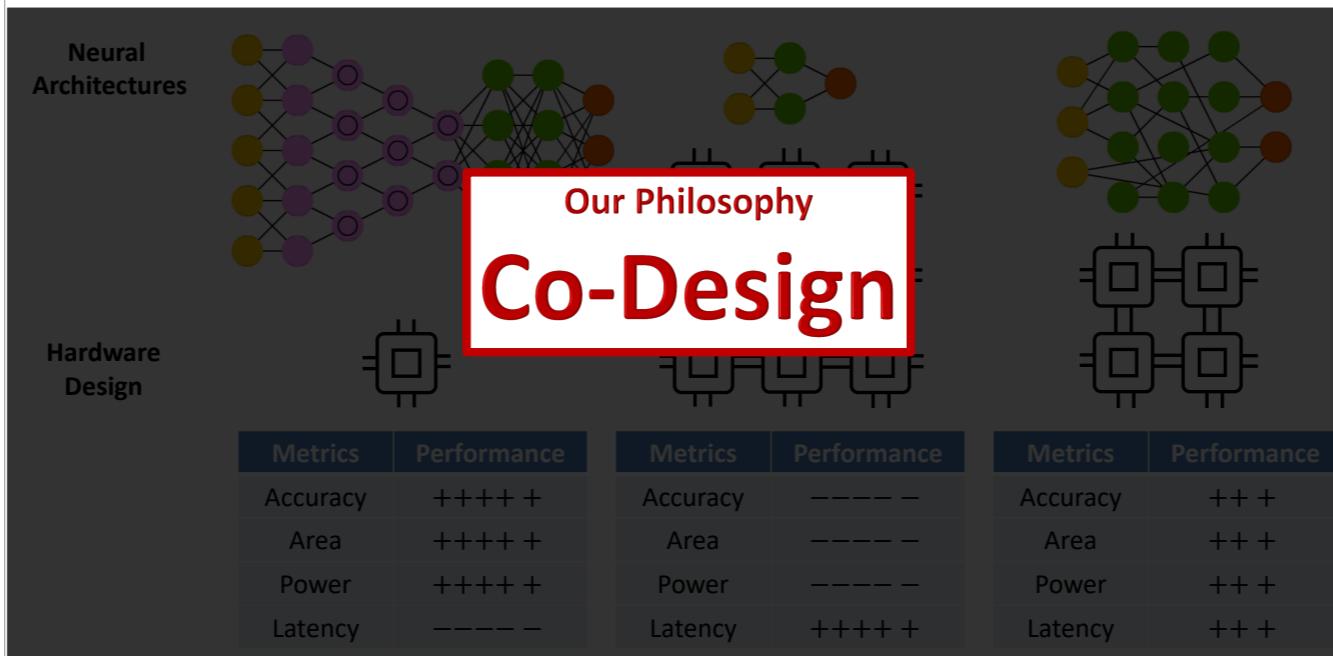
- WE FOCUS ON EDGE AI



8

Recent work for edge computing is extremely restrictive for batter power, area, and accuracy and low latency. Also some will require high security (e.g. healthcare). Trade-offs - low latency and small area, but increased accuracy issues.

- TRADEOFFS BETWEEN ACCURACY AND HARDWARE EFFICIENCY



Trade-offs - if implement left with higher accuracy, power is lower, however non-implementation latency. Next case, smaller network, multiple chips (more hw resources to use) shorter latency, but lower accuracy. Which is worse?

Compromise with third - most practical is to design a neural network, can achieve an acceptable accuracy and also area, power and latency restraints.

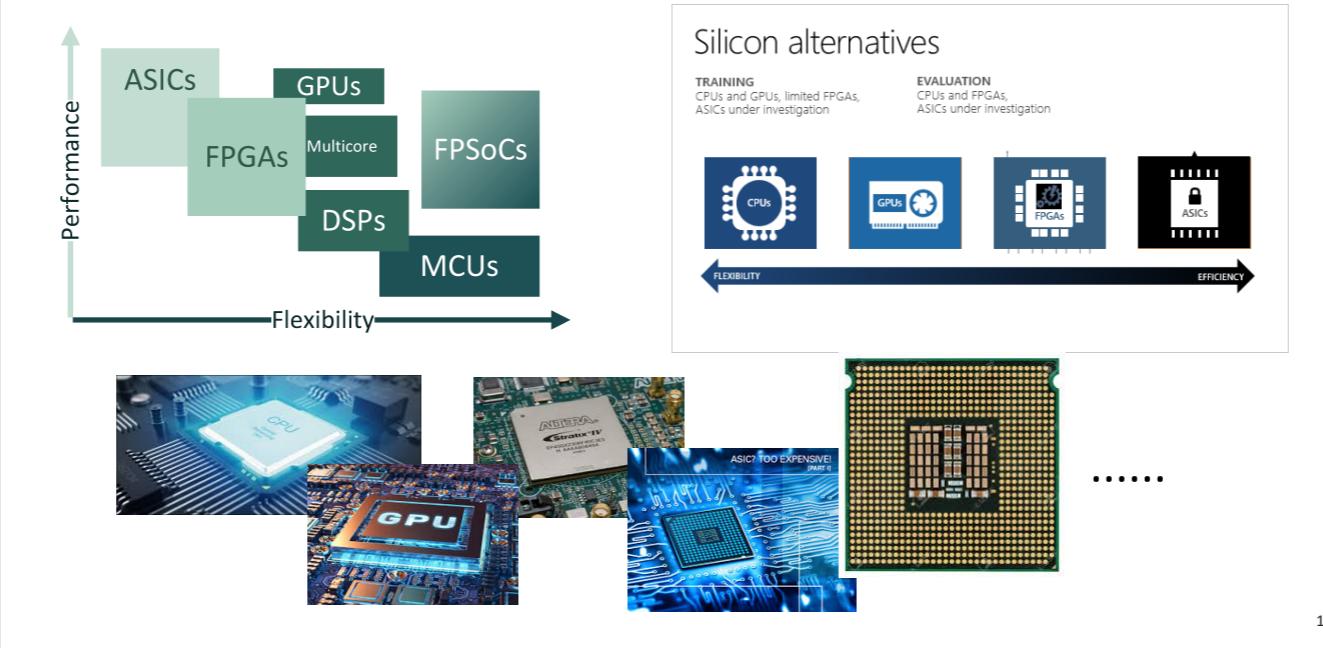
## OUTLINE

- Hardware Design Space 
- Software Design Space
- First Step: Co-Exploration of Neural Architecture and FPGAs
- Second Step: Co-Exploration of Neural Architecture and Heterogenous ASICs
- Third Step: Co-Exploration of Neural Architecture and Scalable ASICs
- Future Work

10

Several frameworks and several steps to solve this in co-exploration.

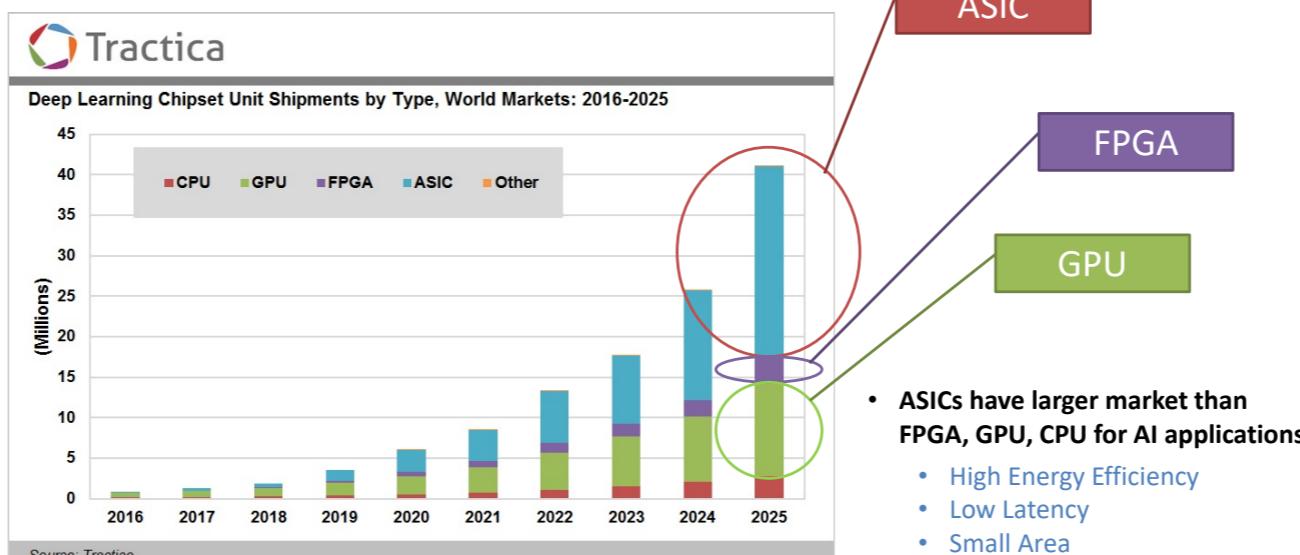
- COMPLEX DESIGN SPACE WITH KINDS OF COMPUTING SYSTEMS



There are existing hardware platforms that can provide variations in power and area constraints. These different computing components execute the

The cpus and gpus are under investigation as as fpgas for implementations

## ■ ASIC WILL HAVE THE LARGEST MARKET



12

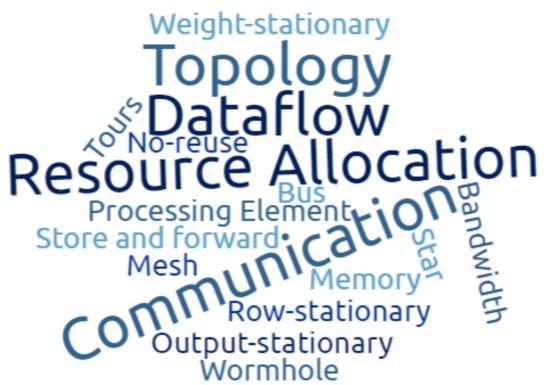
ASICs are composed of processing elements and are connected. Can offer comparable power efficiencies and as a result are becoming more popular. ASICs will dominate the market due to energy, low latency and small area.

Design Space of FPGAs, ASICs and multicores is

# Huge

**Dataflow (data reuse):**

- Weight stationary
- Output stationary
- Row stationary
- No reuse



**Communication:**

- Wormhole
- Store and forward

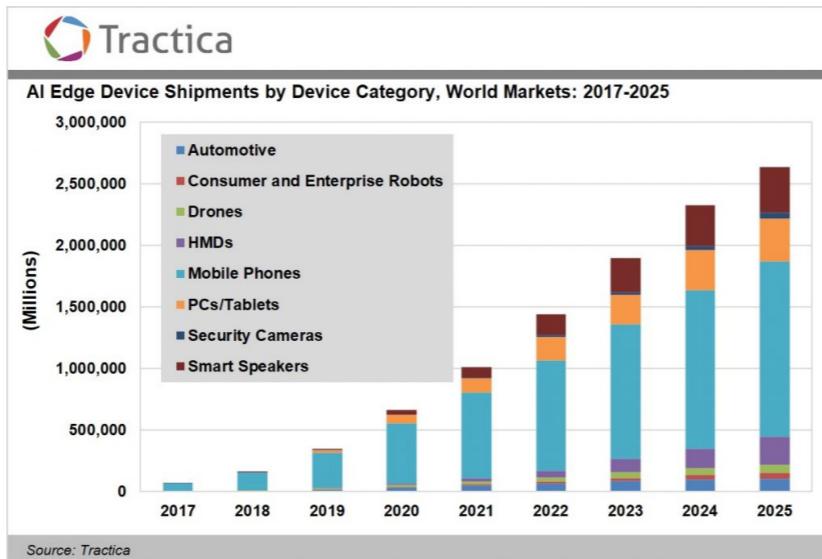
13

Design space of ASICs is huge, have common layer of features. Need to determine the data flow, the memory allocation and topology, as well as the communication between them. This is a challenging task. When ASICs are deployed in edge devices, need to account for various applications being triggered simultaneously (e.g. AI glasses trigger multiple sensors). One data flow cannot fit all tasks. Will require task-level parallelism and heterogenous and multiple accelerators. Requires huge design space on top of just ASIC hardware design space.

## OUTLINE

- Hardware Design Space
- Software Design Space 
- First Step: Co-Exploration of Neural Architecture and FPGAs
- Second Step: Co-Exploration of Neural Architecture and Heterogenous ASICs
- Third Step: Co-Exploration of Neural Architecture and Scalable ASICs
- Future Work

## ■ DEDICATED NEURAL NETWORKS FOR CATEGORIES OF APPLICATIONS



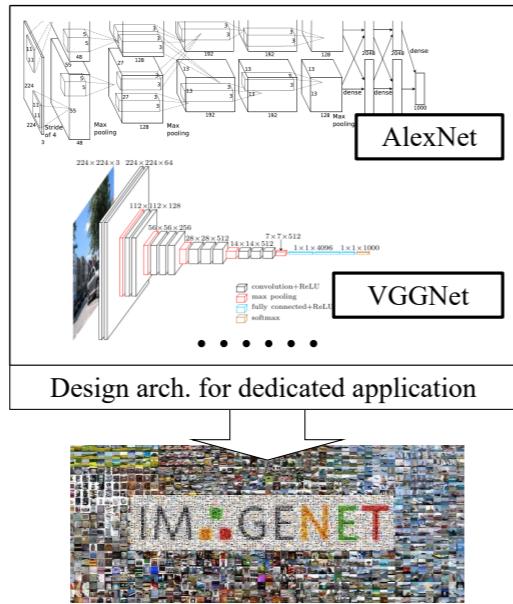
- Increased Categories of edge AI applications
- Applications need **dedicated neural networks** to adapt input data (e.g., Image, Voice)
- Design dedicated neural networks for specific apps?
- To achieve expected performance?



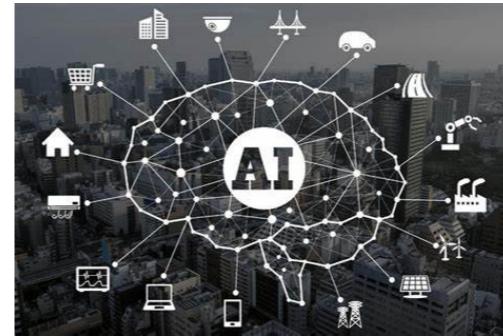
15

# of applications categories is increasing (e.g. voice and image require different software). How to design specific neural network to achieve. Neural architecture search

## ■ HUMAN INVENTED NEURAL ARCHITECTURES



Design arch. for dedicated application



Era of AI Democratization

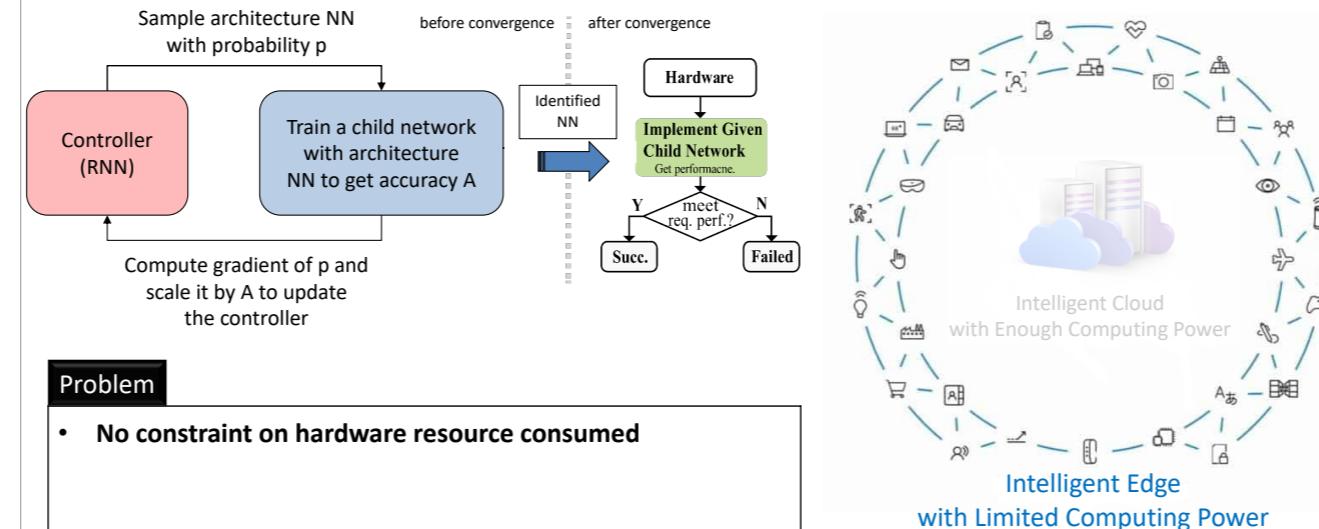
### Problem

- Domain knowledge and excessive labor
- It is impossible to manually design specific arch. for each dedicated application in the era of AI democratization

16

Neural networks designed for image classifications.

## ■ NEURAL ARCHITECTURE SEARCH (NAS)

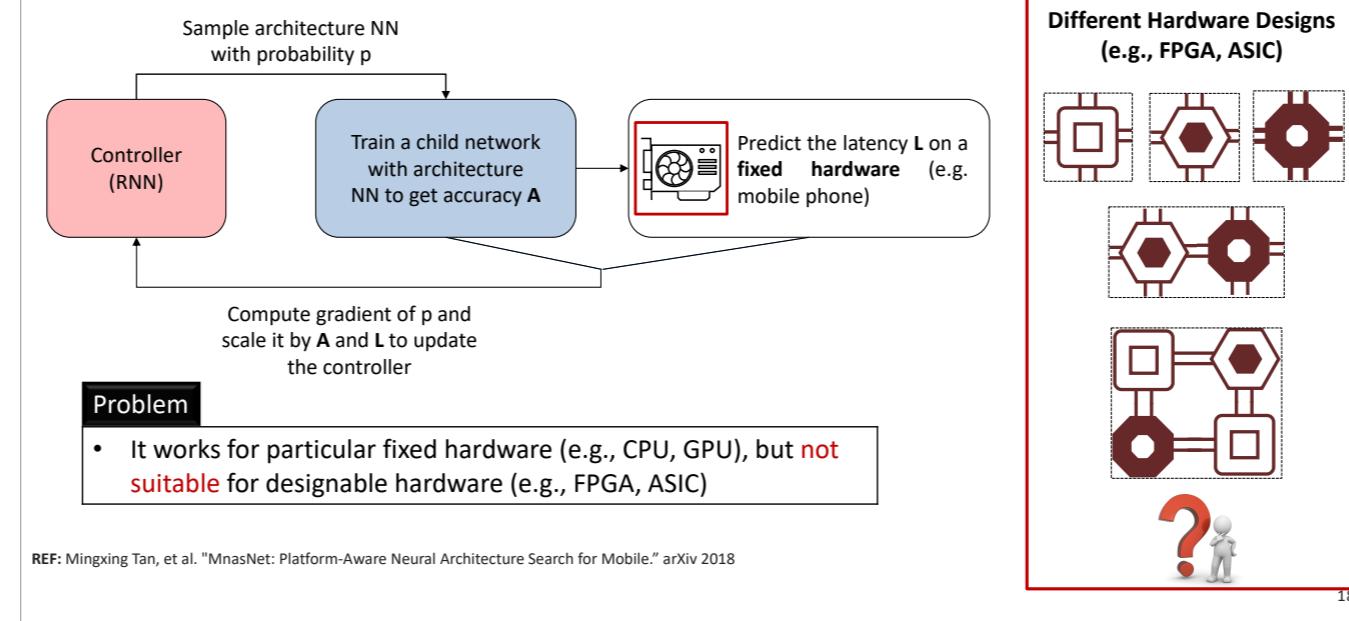


REF: Zoph, Barret, and Quoc V. Le. "Neural architecture search with reinforcement learning." *ICLR 2017*

17

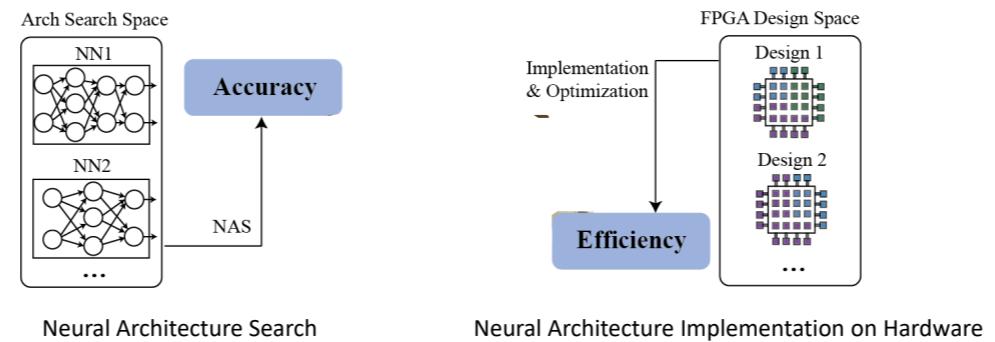
Recently, instead of human-designed neural nets, architecture search that would allow for automatic approach. Hyperparameters train a child network, accuracy will be fed back into a reinforce learning. Fine for cloud computing with plenty of power, but with edge computing, it is limited (latency and throughput). Moving NAS to edge devices violates hw specifications.

## ■ HARDWARE-AWARE NAS



Further approach - hardware-aware NAS. Hardware efficiency is considered in hw design. Power results will also be sent back to controller to further refine neural network. Caveat - designed for fixed hardware platforms. A simple model for latency cannot determine with scenario.

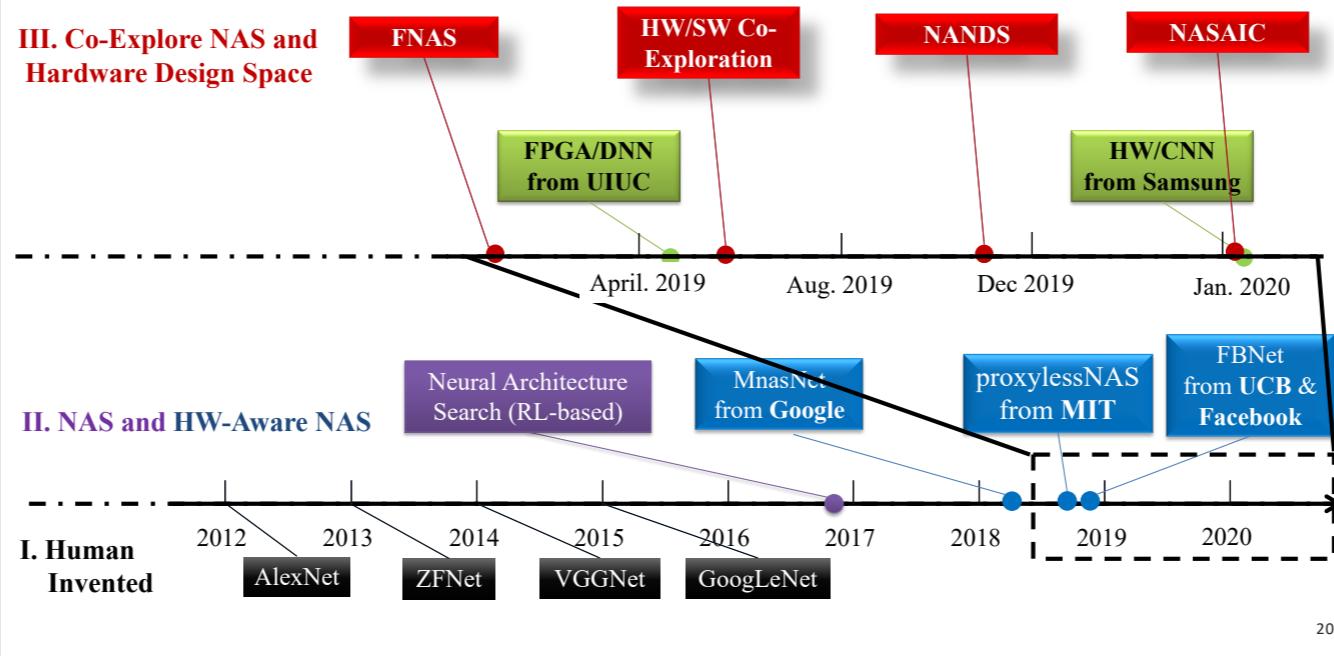
## ■ A MISSING LINK BETWEEN TWO DESIGN SPACES



19

Software-based approach can determine best neural network to use, but combining with the flexible hardware is missing.

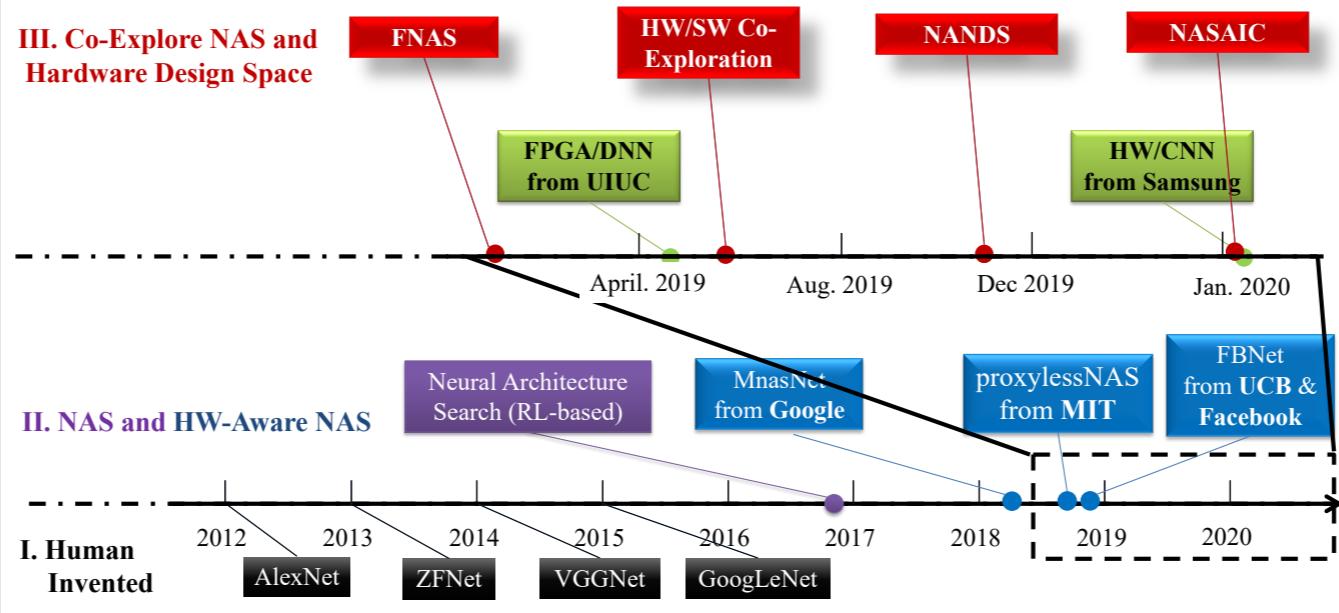
## ■ EVOLUTION OF DEEP NEURAL ARCHITECTURE EXPLORATION



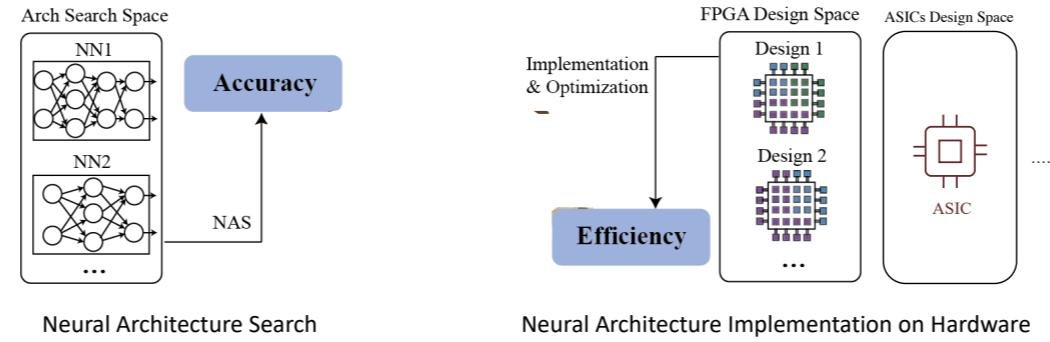
Summary:

- 1<sup>st</sup> stage is human-invented stage: neural net is designed based on human
- 2<sup>nd</sup> NAS approach (circa 2017) reinforcement learning in controller loop (fixed hw only)
- 3<sup>rd</sup> Co-explore of NAS and hw design simultaneously (e.g. network on chip)

## ■ EVOLUTION OF DEEP NEURAL ARCHITECTURE EXPLORATION



## ■ A MISSING LINK BETWEEN TWO DESIGN SPACES



## OUTLINE

- Hardware Design Space
- Software Design Space
- First Step: Co-Exploration of Neural Architecture and FPGAs 
- Second Step: Co-Exploration of Neural Architecture and Heterogenous ASICs
- Third Step: Co-Exploration of Neural Architecture and Scalable ASICs
- Future Work

## ■ FPGAs IN DNN APPLICATIONS

FPGA in Cloud Computing



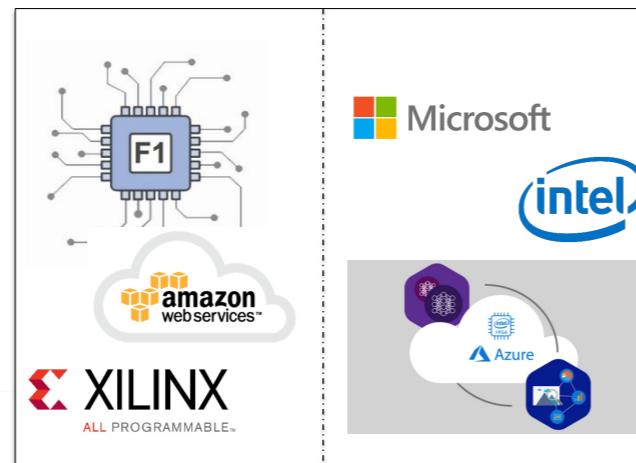
[ref] Where Do FPGAs Stand in Auto IC Race? [https://www.eetimes.com/document.asp?doc\\_id=1333419#](https://www.eetimes.com/document.asp?doc_id=1333419#)  
[ref] PYNQ in UAV. <http://brennancain.com/pynqcopter-an-open-source-fpga-overlay-for-uavs/>

25

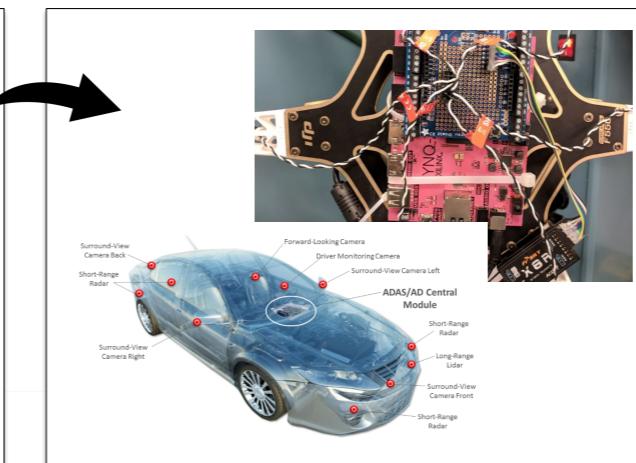
Can be configured by customer, specified by HDL. Relevant to cloud computing and big data. Pressure on the amount of data (rapid growth for computation and calculation) especially with AI applications. Should have high throughput (capacity) to process large amounts of data and also low latency to respond to network device immediately, and have low power for high density cloud computing. Performancewise FPGA is not restricted by Von Neumann and has high performance of ASIC, but with flexibility. Suited for data-level parallelism. Allow for multiple row computing and adapts to the evolving networking and storage algorithms. Optimized algorithms offer better power to computation ratio. Flexibility is area where FPGAs are most dominant over ASICs. All of these mean FPGAs are pathway to accelerating cloud computing (e.g. in 2004 Microsoft developed 1000 servers for search based on FPGAs increasing throughput by 90%). FPGAs are widely used in edge computing as well (not as fast, but flexible).

## ■ FPGAs IN DNN APPLICATIONS

### FPGA in Cloud Computing



### FPGA in Edge Computing

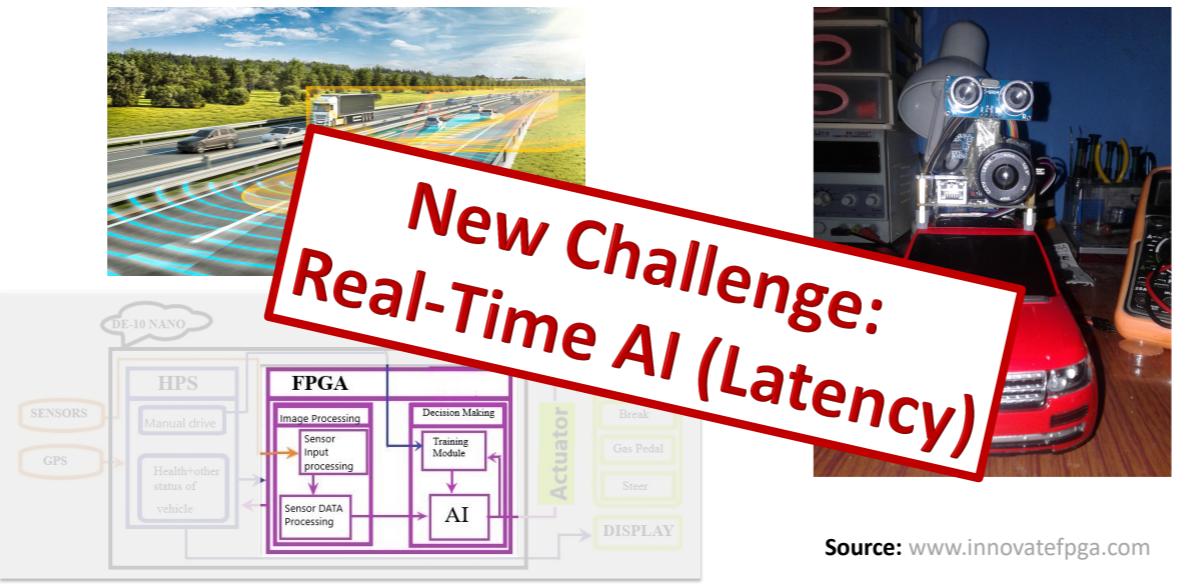


26

[ref] Where Do FPGAs Stand in Auto IC Race? [https://www.eetimes.com/document.asp?doc\\_id=1333419#](https://www.eetimes.com/document.asp?doc_id=1333419#)  
[ref] PYNQ in UAV. <http://brennancain.com/pynqcopter-an-open-source-fpga-overlay-for-uavs/>

FPGAs are suitable for edge computing and IoT including self-driving. Cloud-based solutions will not be able to provide time-sensitive operations. Consensus is to expand cloud services to be deployed to the edge. Fog or edge computing. Current latency intensive, (critical edge servers) requires GPUs, but cannot deliver predictable results (specifically they are power hungry).

- FPGAs IN EDGE AI: AUTONOMOUS CAR

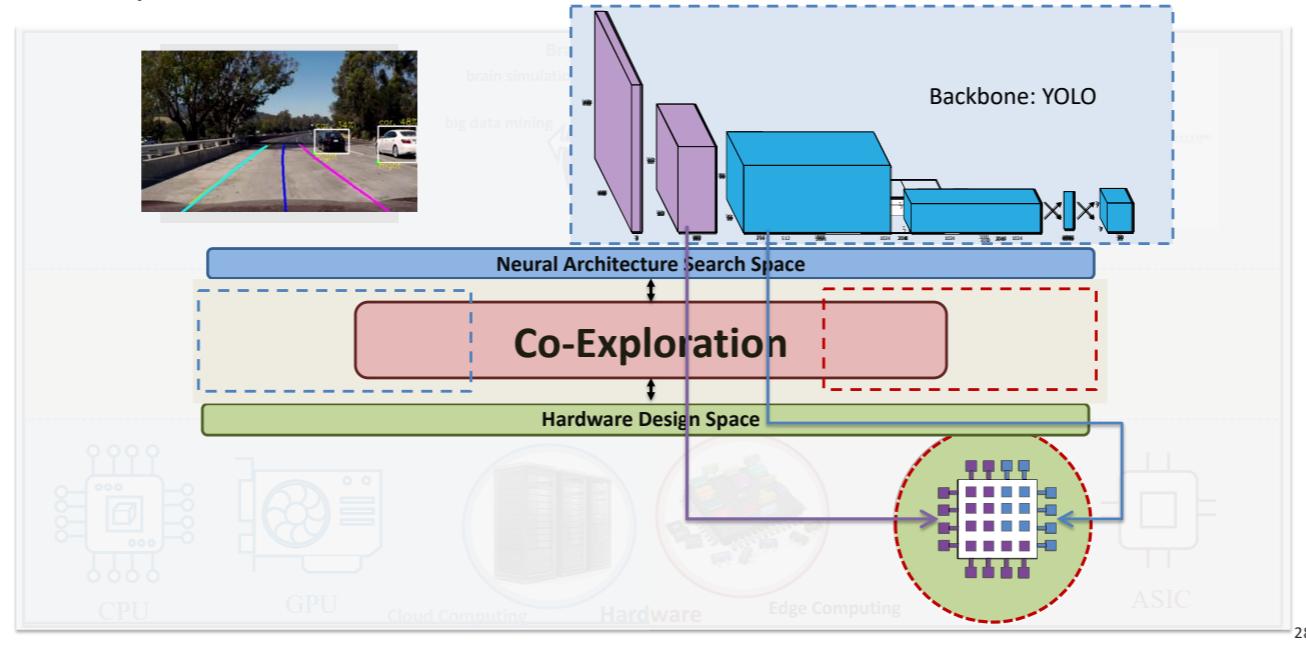


Source: [www.innovatefpga.com](http://www.innovatefpga.com)

27

FPGAs are widely used in autonomous cars. AI chip manufacturers have realized the latency issues of standard chips. Currently rely on data centers far away; adds dangerous wait, but also enormous power constraints. Need to address critical mission tasks accurately and quickly.

## ■ HW/SW Co-EXPLORATION

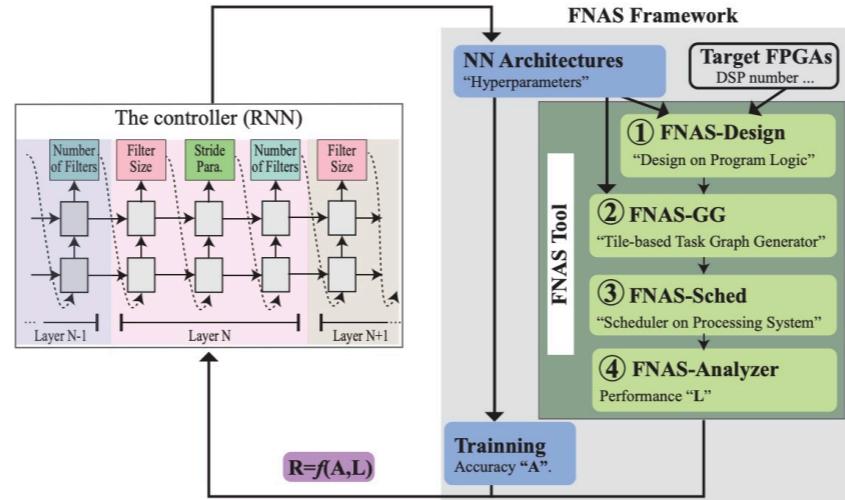


Co-exploration: FPGA is the platform to implement the NAS network. Include additional matrix to describe the neural net, but can be challenging. Unlike human-developed, NAS can be irregular and design flows that are not suitable for FPGA (large sizes requiring multiple FPGAs). Need an abstraction model to achieve efficient model (hardware optimization loop tiling). Dependancies can be captured by extra fabric. NAS can be parallelized, but FPGA doesn't realize that. Analyze pipeline stores and scheduling to estimate the latency. Generate optimal architecture and optimal FPGAs. In the framework, it can model the different kinds of tasks.

Hardware-aware NAS using FPGA as a platform (high performance and energy efficiency).

Include additional matrix for neural architecture search. Needs multiple FPGA necessary and therefore the need to coordinate these. Need a bridge between the NAS search and between the hardware. Starting the design principle of schedulers for scheduling tasks and estimate overall latency. Can also model other types of platforms.

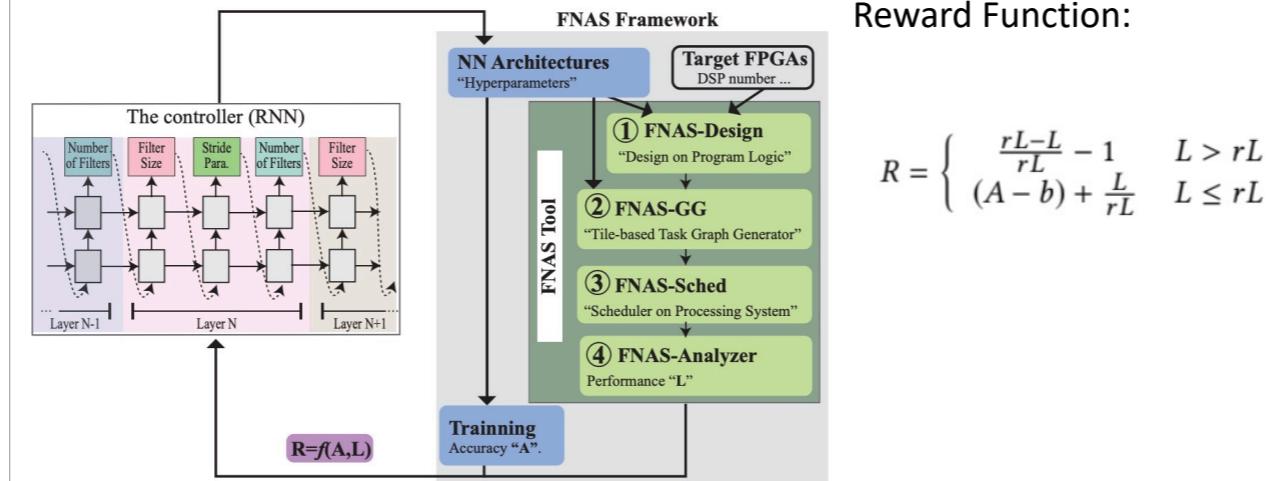
## ■ CO-EXPLORATION FRAMEWORK



29

Develop FPGA with NAS with required latency.  $R$  = reward (based on Accuracy and Latency).

## ■ CO-EXPLORATION FRAMEWORK



Reward Function:

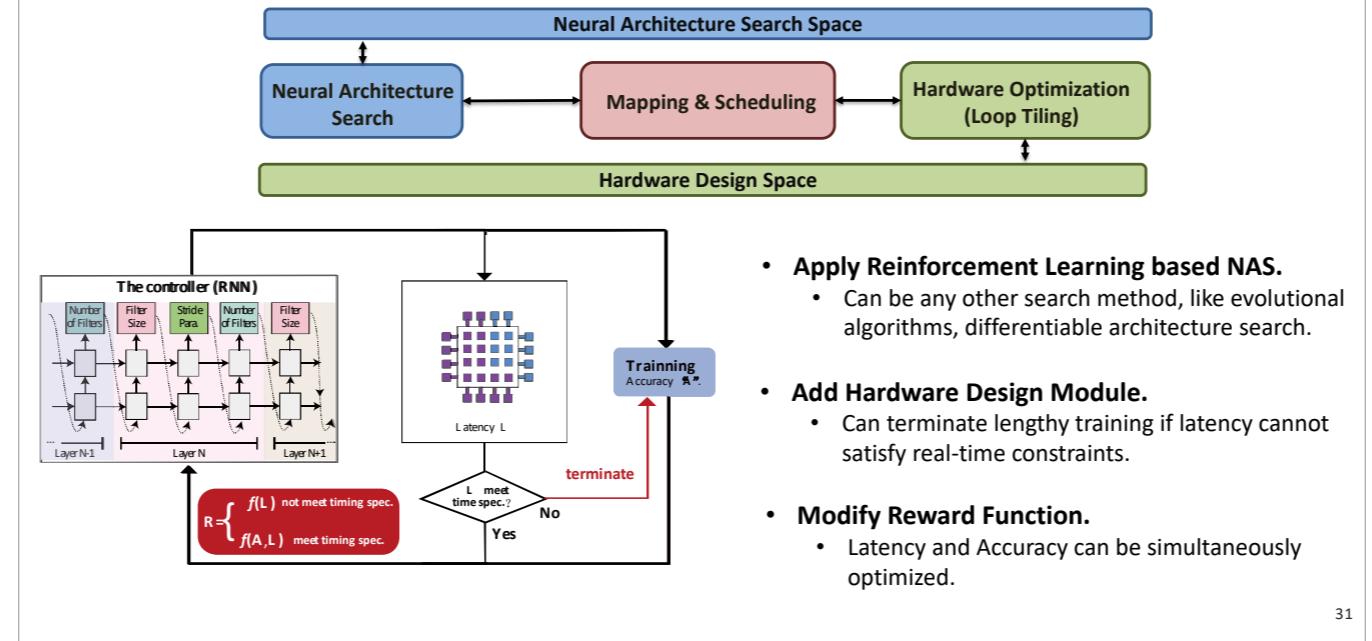
$$R = \begin{cases} \frac{rL-L}{rL} - 1 & L > rL \\ (A - b) + \frac{L}{rL} & L \leq rL \end{cases}$$

30

There are two cases, 1<sup>st</sup> performance of the application [...] return

Solution has higher performance reward if its latency is less than the

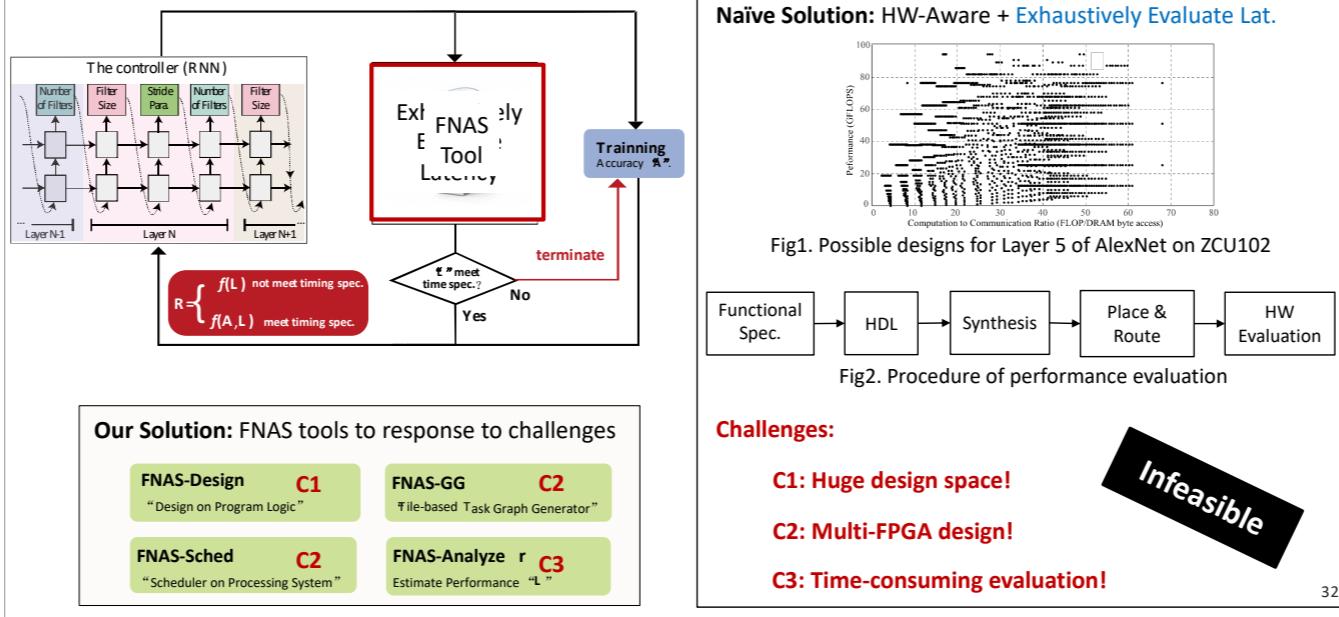
## Co-EXPLORATION FRAMEWORK



In the implementation of the proposed co-explor the reinforcement learning can be (e.g.,

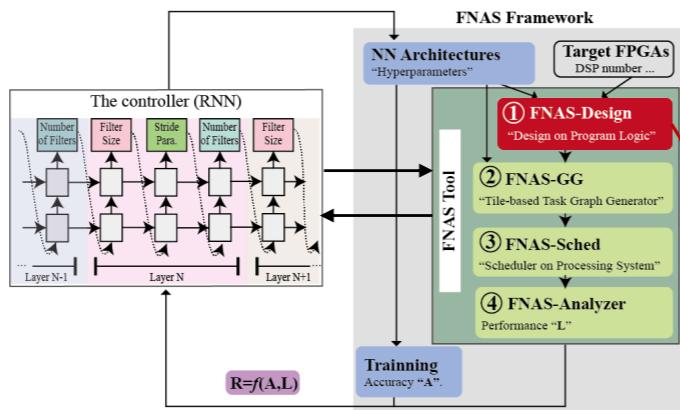
Take different amnt of resources depending on the complexity of the algorithm. New hdwr design modules are added as the

## PERFORMANCE ESTIMATION: SOLUTIONS & CHALLENGES



32

## FNAS: DESIGN OPTIMIZATION



REF: Chen Zhang et al. 2015. Optimizing fpga-based accelerator design for deep convolutional neural networks. In Proc. of FPGA.

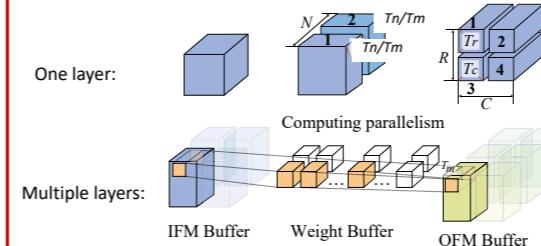
### 1 FNAS Design

Designed to determine the tiling parameters for a given NN architecture on target FPGAs.

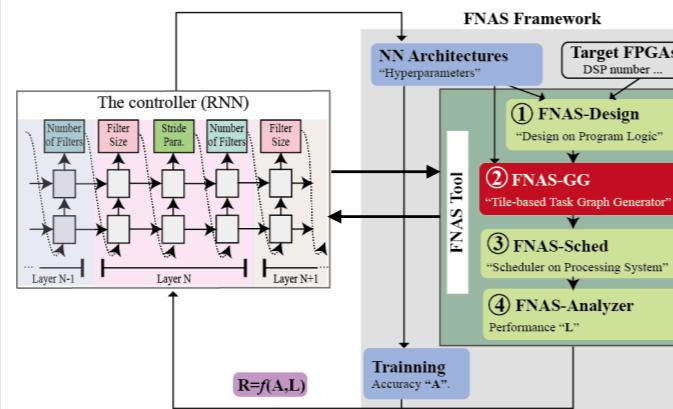
#### On-chip accelerator design:

##### Determine:

1. On-chip buffer allocation; 2. Accelerator size for computing  
(note: both are determined by tiling parameters,  $T_m$ ,  $T_n$ ,  $T_r$ ,  $T_c$ )



## FNAS: GRAPH GENERATOR

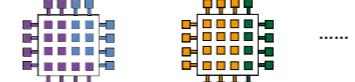


### ② FNAS GG

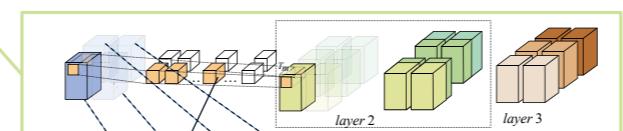
a graph generator that takes the design parameters and NN architecture to generate the dependency graph between data tiles and tasks

### Given:

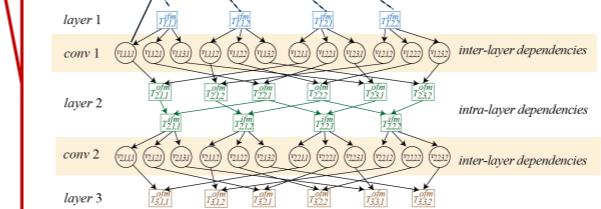
1. FPGAs with attributes including LUTs, DSPs, BRAM, etc.



2. A neural architecture with determined hyperparameters



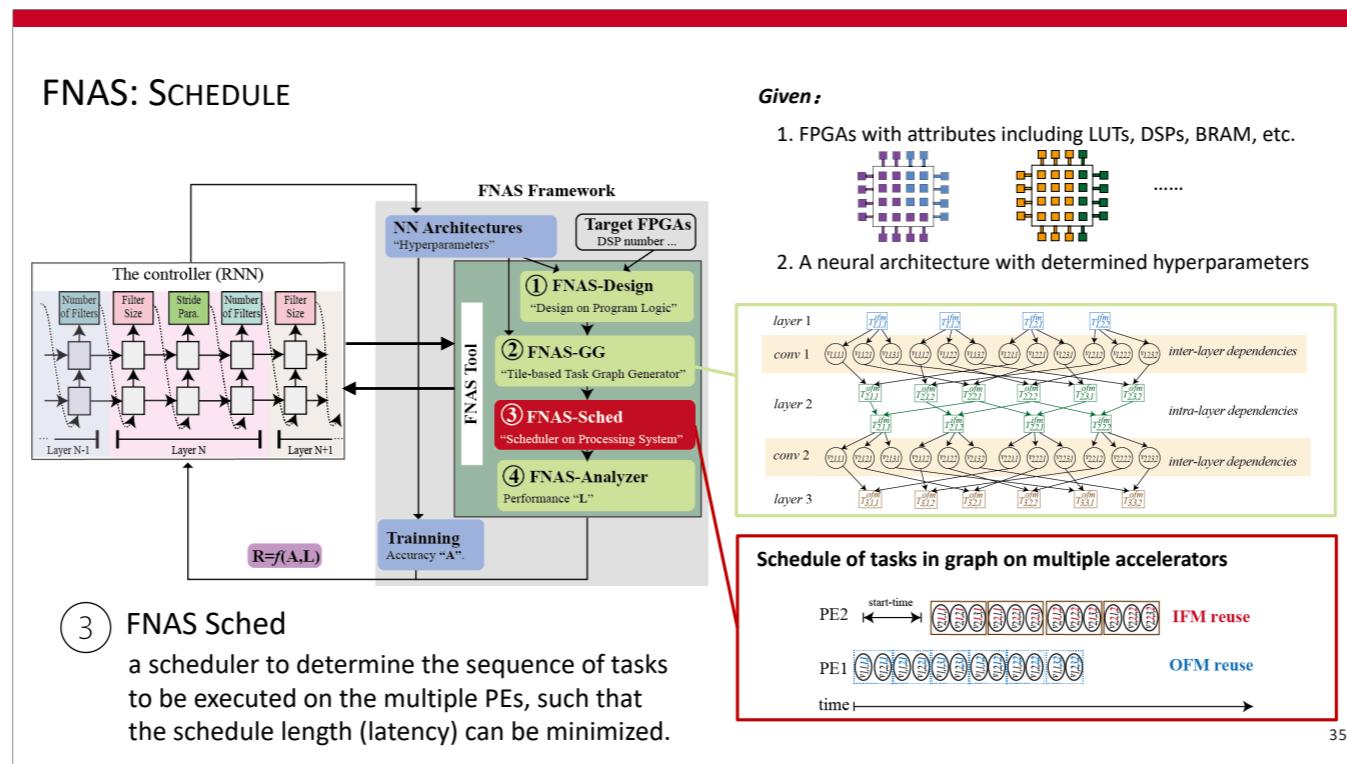
### High-level graph abstraction



34

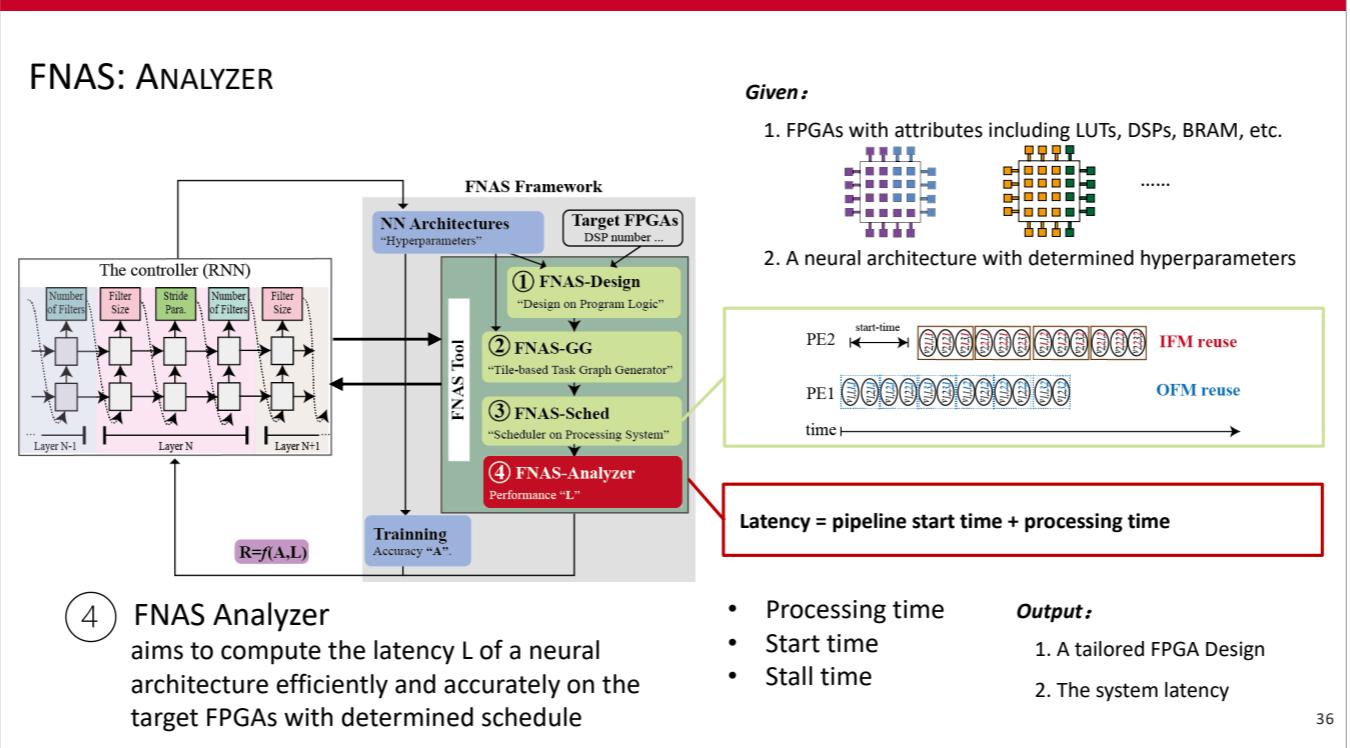
After the tuning, the next component will generate the dependency graph. Two kinds of dependencies: interlayer (in 2 convolutional layers) and intralayer (FNAs tries to maximize parallelism based on the previous input

1. minimize processor start time
2. Maximize the reuse of data
3. Minimize the pipeline stores (by lack of input data on 1 processor)



Two data read strategies 1. Output future map (same output future map tile - uniform reuse across all tiles and maximize data reuse) 2. Input future map (

Accelerator can be implemented for further optimization



Ready-to-run to run tasks that are ready (one task in the queue) to minimize stall time.

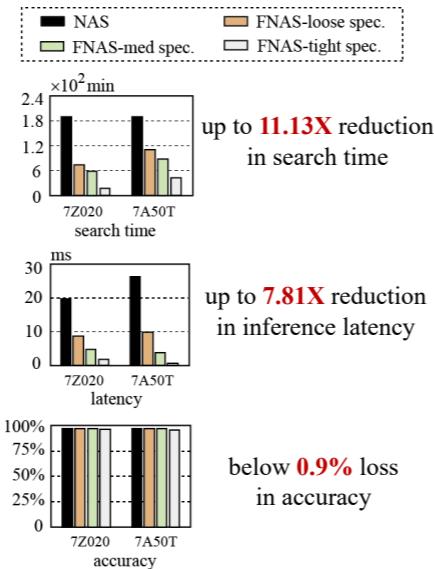
FNAS determines if the latency requirement will be met, if not, then quits.

## EXPERIMENTAL SETTING

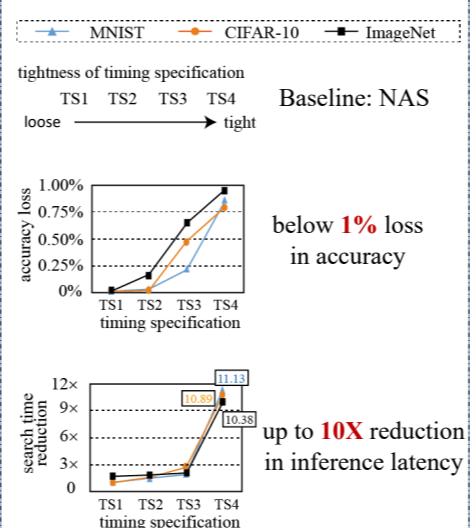
FPGAs		Xilinx 7A50T		Xilinx 7Z020
Datasets				ImageNet
MNIST	up to 5	up to 10	up to 15	
CIFAR-10	[5, 7, 14]	[1, 3, 5, 7]	[1, 3, 5, 7]	
ImageNet	[9, 18, 36]	[24, 36, 48, 64]	[16, 32, 64, 128]	
<b>NAS Search Space</b>	Layer Num.	up to 5	up to 10	up to 15
	Filter Size	[5, 7, 14]	[1, 3, 5, 7]	[1, 3, 5, 7]
	Filter Num.	[9, 18, 36]	[24, 36, 48, 64]	[16, 32, 64, 128]
<b>HW Search Space</b>	Channel Tiling Para. (Tm,Tn); Row Tiling Para. (Tr); Col Tiling Para. (Tc); Schedule			
<b>Timing Spec. (ms)</b>	[2, 5, 10, 20]	[1.5, 2, 2.5, 10]	[2.5, 5, 7.5, 10]	37

## EXPERIMENTAL RESULTS

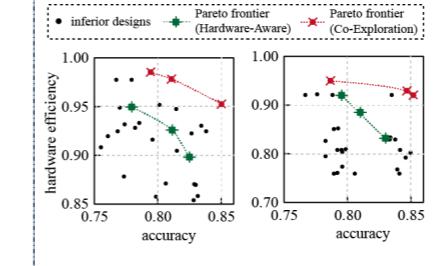
### Different Hardware (MNIST)



### Different Datasets (7Z020)



### Compare to HW-Aware NAS (CIFAR-10 + 7Z020)



38

Two reasons for improvement: Early stage pruning (terminate training if latency will not be met) and

## EXPERIMENTAL RESULTS: SUPERIOR TO EXISTING APPROACHES

### Optimizing Hardware Efficiency

Comparison the proposed Co-Exploration with Hardware-Aware NAS and Heuristic Sequential Optimization

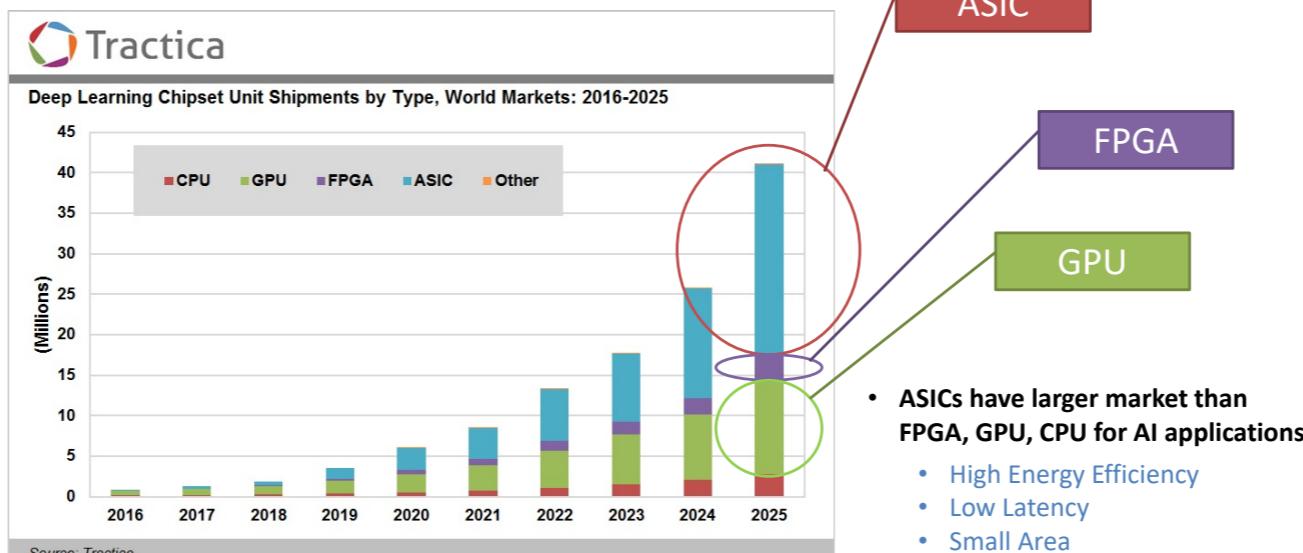
Dataset	Models	Depth	Parameters	Accuracy (Top1)	Accuracy (Top5)	Pipeline Eff.	FPS	Energy Eff. GOPs/W
CIFAR-10	Hardware-Aware NAS	13	0.53M	84.53%	-	73.27%	16.2	0.84
	Sequential Optimization	13	0.53M	84.53%	-	92.20%	29.7	1.36
	Co-Exploration (OptHW)	10	0.29M	80.18%	-	99.69%	35.5	2.55
	Co-Exploration (OptSW)	14	0.61M	85.19%	-	92.15%	35.5	1.91
ImageNet	Hardware-Aware NAS	15	0.44M	68.40%	89.84%	81.07%	6.8	0.34
	Sequential Optimization	15	0.44M	68.40%	89.84%	86.75%	10.4	0.46
	Co-Exploration (OptHW)	17	0.54M	68.00%	89.60%	96.15%	12.1	1.01
	Co-Exploration (OptSW)	15	0.48M	70.24%	90.53%	93.89%	10.5	0.74

### Optimizing Network Accuracy

## OUTLINE

- Hardware Design Space
- Software Design Space
- First Step: Co-Exploration of Neural Architecture and FPGAs
- Second Step: Co-Exploration of Neural Architecture and Heterogenous ASICs 
- Third Step: Co-Exploration of Neural Architecture and Scalable ASICs
- Future Work

- ASIC WILL HAVE THE LARGEST MARKET



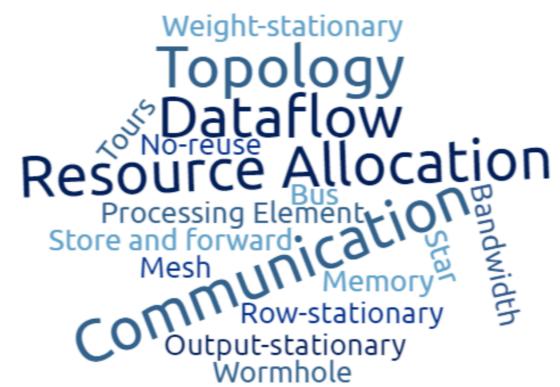
43

ASICs provide incomparable energy efficiency, latency, and design area. Should push toward ASIC to accommodate them.

# Design Space of ASIC is Huge

## Dataflow (data reuse):

- Weight stationary
- Output stationary
- Row stationary
- No reuse



## Communication:

- Wormhole
- Store and forward

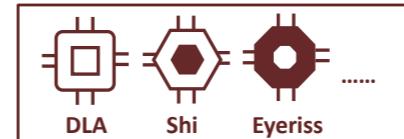
44

Most ASIC work do not take into account data flow, topology between processing elements, and data flow (and data reuse). ASIC accelerators deployed at edge need to be able to handle multiple tasks (such as image detection) concurrently. Requires task-level parallelism require multiple accelerators. Trying to combine these accelerators into one for energy efficiency has been difficult.

## ▪ MOTIVATION: TEMPLATE POOL

Existing ASIC Accelerators for Neural Network			
Dataflow ➔	weight stationary (WS)	output stationary (OS)	row stationary (RS)
Template ➔	DLA (NVIDIA'17) nn-X, TPU, .....	ShiDianNao (CAS&EPFL'15) Gupta, Peemen, .....	Eyeriss (MIT'16) Eyeriss v2, .....
			DianNao (CAS'14) DaDianNao, Zhang, .....

 Well-designed ASICs → No need to design ASIC from scratch!  
Just select **the ASIC design** (called Template) from the pool



Template Pool

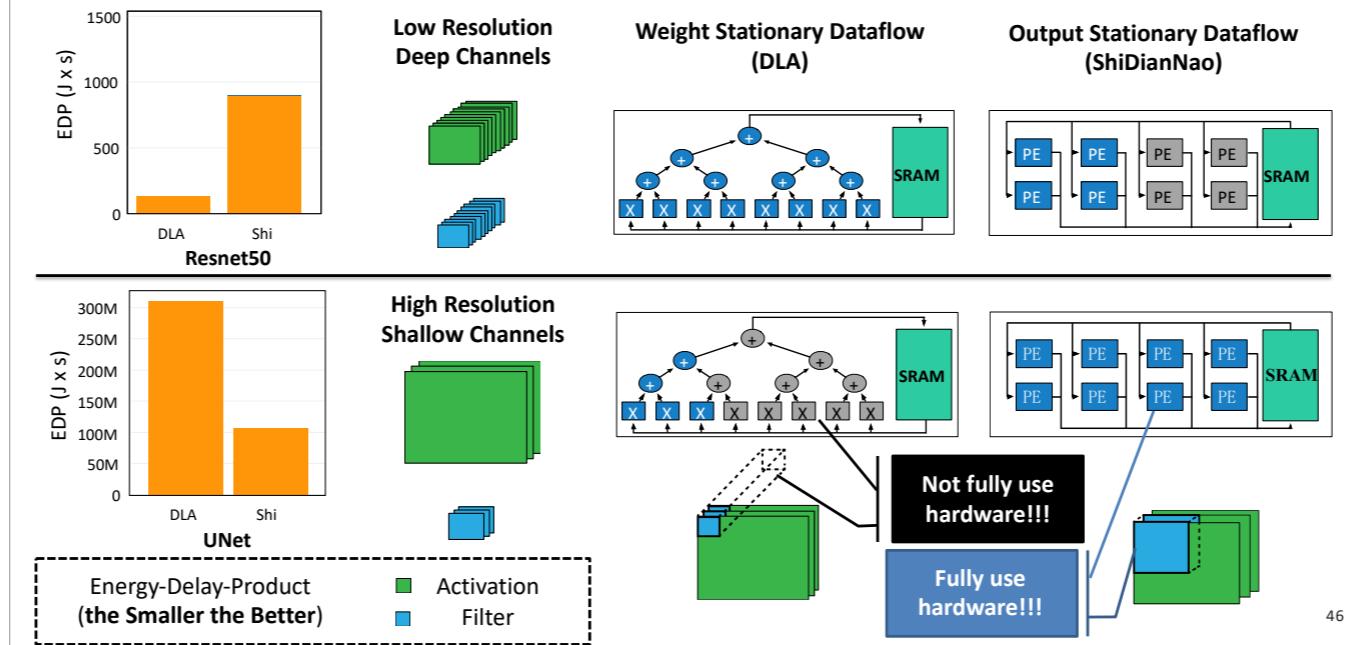


**One template for all neural network architectures?**

45

Use existing ASIC data flows to build a pool of ASIC templates.

## MOTIVATION: HETEROGENEOUS ASICs



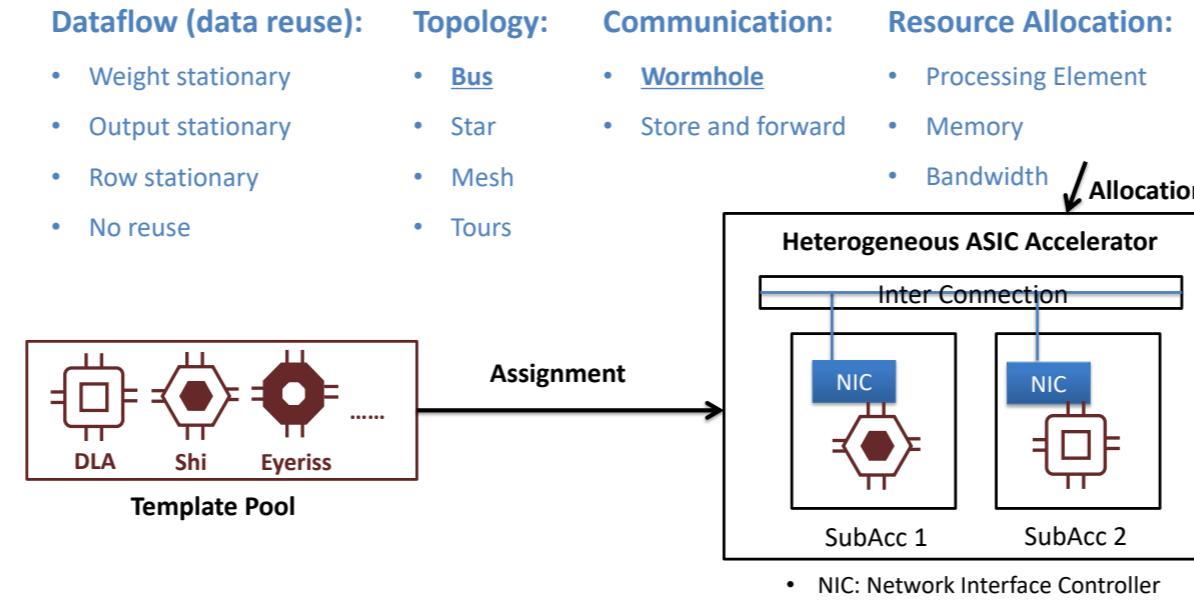
46

Low resolution and deep channels

ShiDianNao is high resolution and shallow channels

Due to distinct futures, variety of methods should be deployed onto hwre platforms, fully using hardware and handling multiple tasks.

- Design Space: Take ASIC Accelerator with 2 Sub-Accelerator as An Example



Design space of ASICs is huge. Interconnect two ASIC accelerators with pool of ASIC templates, we need to allocate the memory and compute resources (data bandwidth and memory size).

## ■ PROBLEM AND CHALLENGES

Realistic Problem from  
On-Device AI Group  
(AR/VR Glasses)



### Given Design Space:

- Datasets for multiple machine learning tasks
- A set of PEs and total Bandwidth

### Given Constraints:

- Latency; Power; Area.



### Output:

- Neural architecture with the maximum accuracy
- ASIC chip design to satisfy hardware specifications

### Challenge1: ASIC has huge design space.

For the same # of PEs, it can have

- Different topologies (Bus, NoC...)
- Different dataflows (WS, OS...)

### Challenge2: Multiple tasks in application

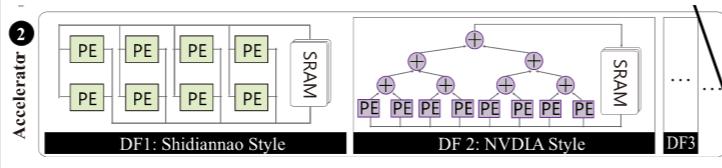
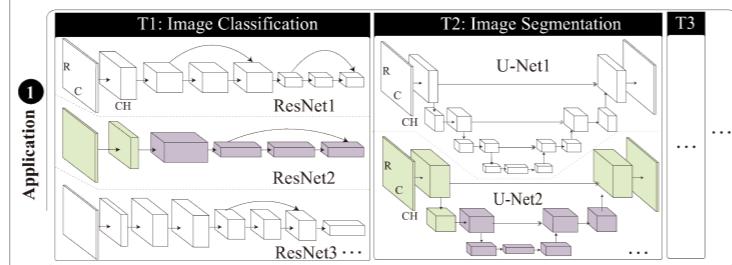
- They subject to the unified constraints (latency, power, area)



### Existing NAS:

**NEITHER ASICs NOR multi-tasks**

## ■ PUT ALL TOGETHER

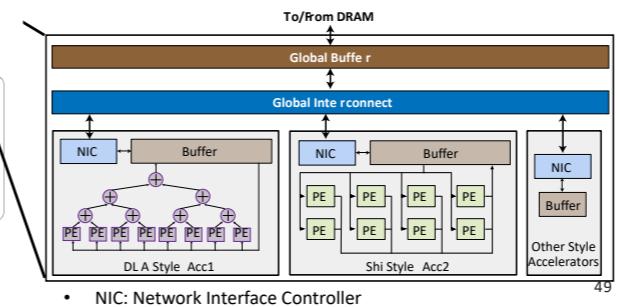


### Given:

- Multiple machine learning tasks
- A pool of ASIC templates
- Unified Constraints: Latency; Power; Area.

### To:

- Determine neural arch. for each task (NAS)



• NIC: Network Interface Controller

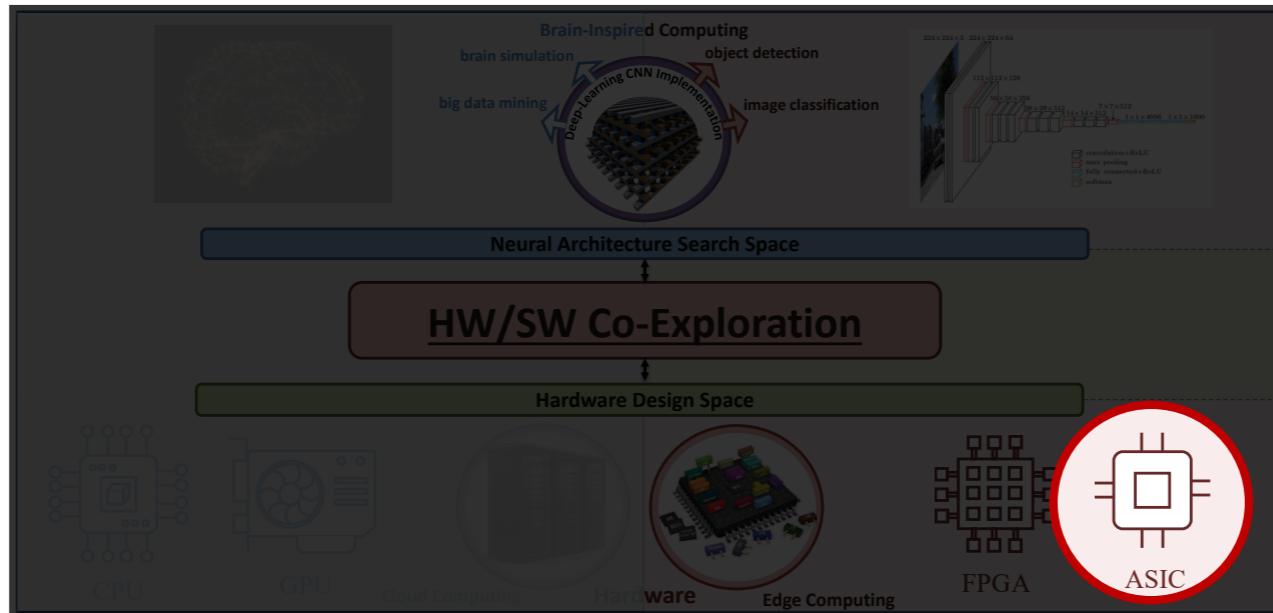
Three exploration layers:

Application: Determine NAS for each task

Accelerator: Creates ASIC template based on existing templates to each sub-accelerator and map to scheduler

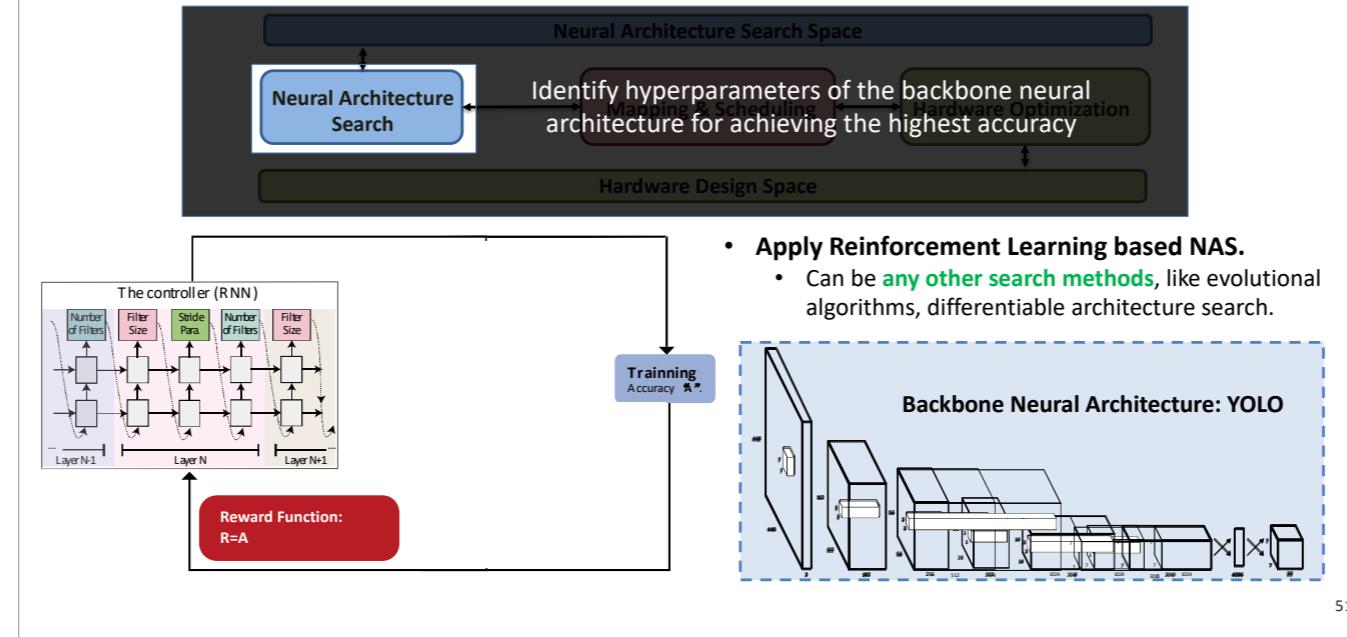
- : resource allocation to templates (PEs bandwidth)
- : mapping and scheduling

## HW/SW Co-EXPLORATION TARGET ASICs



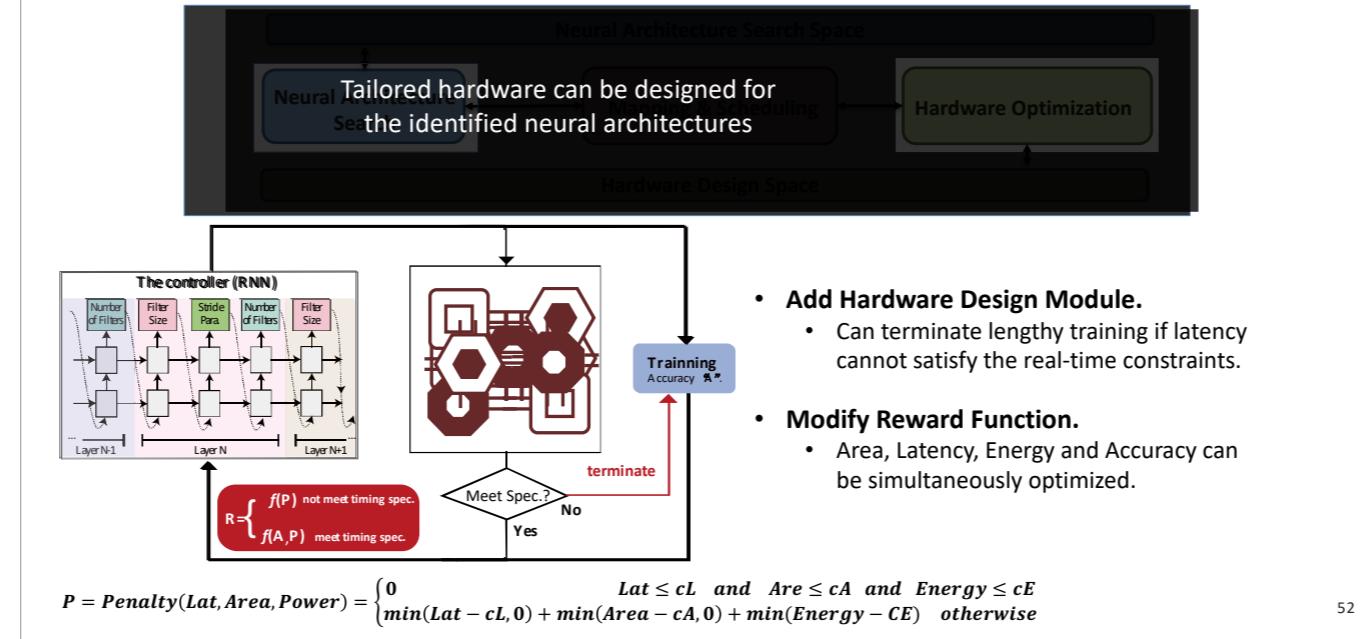
50

## ■ HW/SW Co-EXPLORATION: NAS



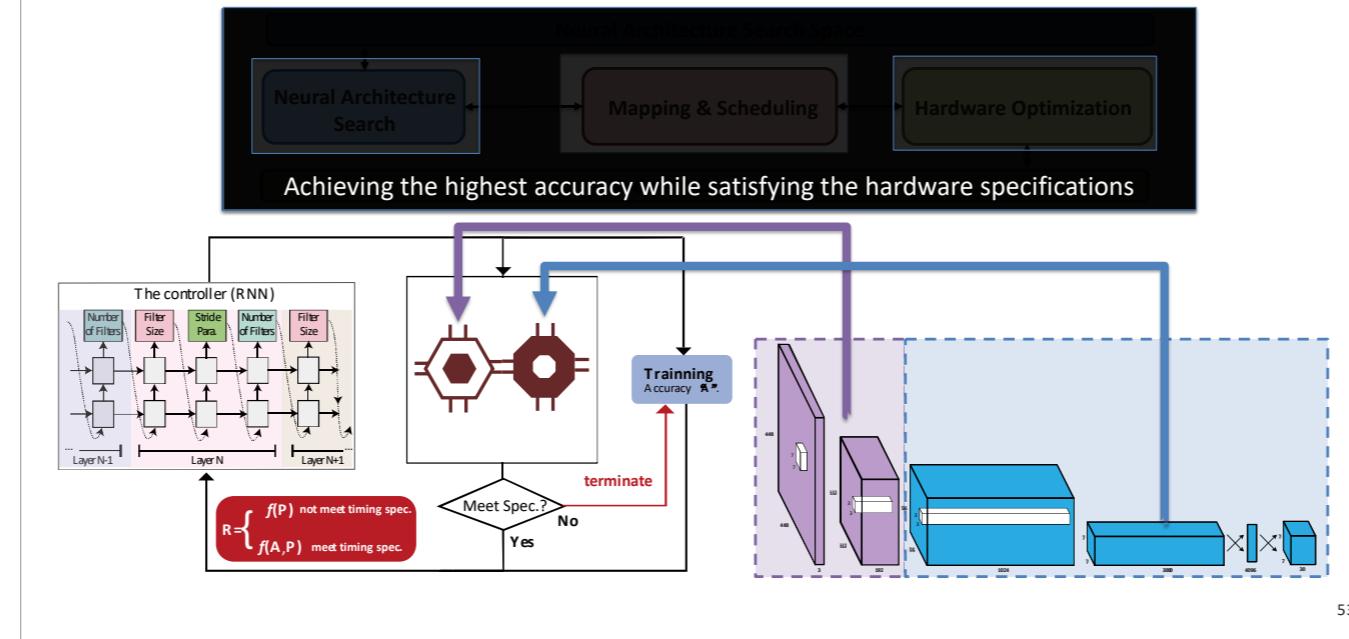
For NAS implemented reinforcement learning and identifying hyperparameters for highest accuracy

## ■ HW/SW CO-EXPLORATION: HW OPTIMIZATION



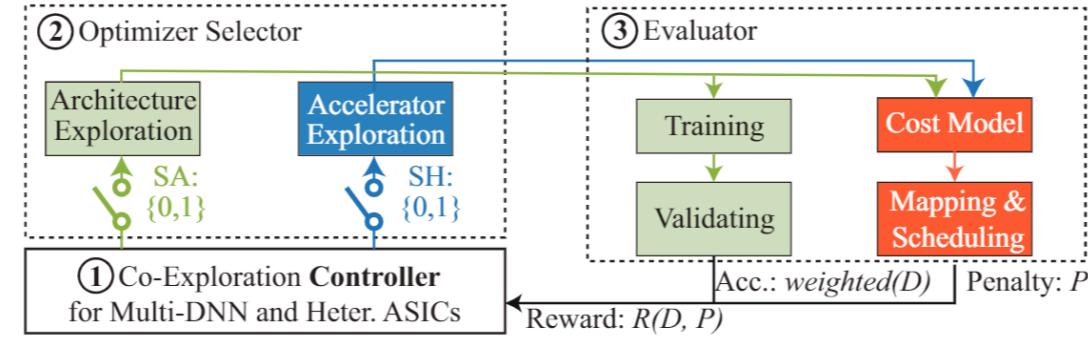
Reduce latency or terminate. Latency, area, and Energy all used to determine the reward

## ■ HW/SW CO-EXPLORATION: MAPPING AND SCHEDULING



53

## ■ PUT ALL TOGETHER: NOVEL NASAIC FRAMEWORK



❖ Controller: sample NN and allocate hardware resource in each iteration

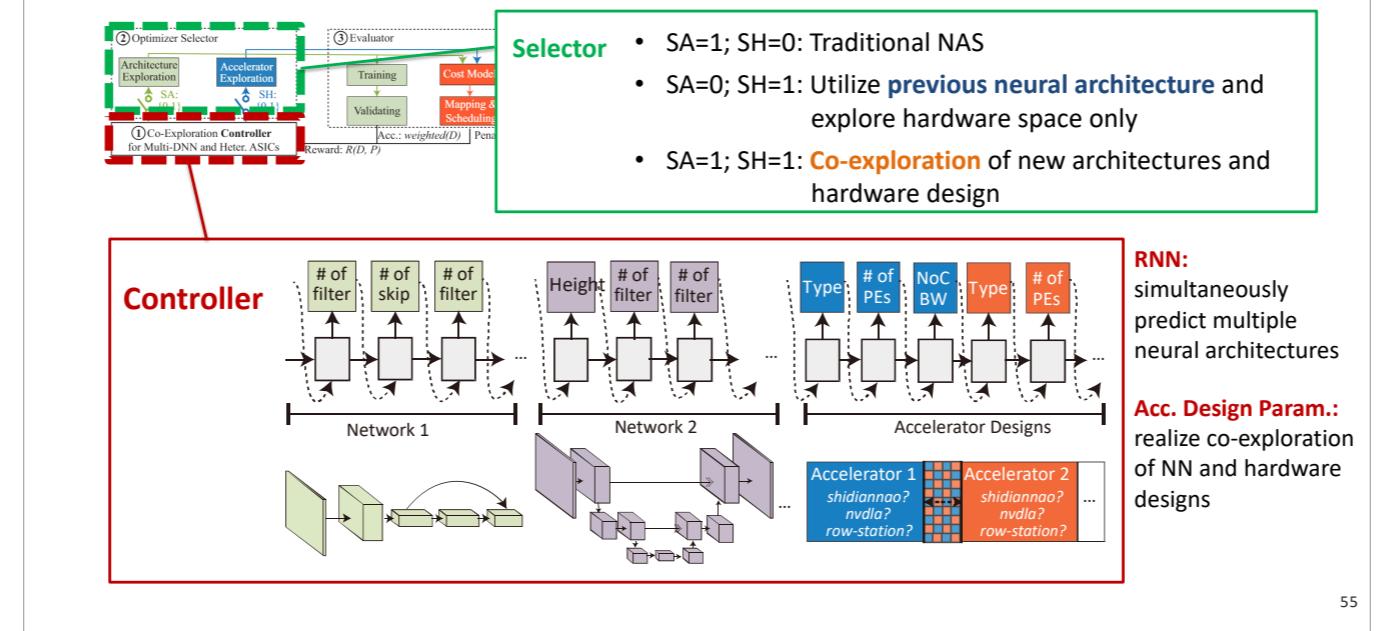
❖ Optimizer selector: NAS and hardware optimization

❖ Evaluator: generate the accuracy and hardware cost

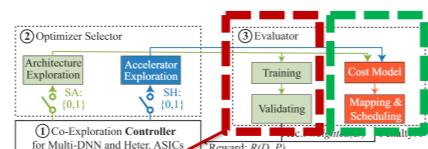
❖ Finally, a reward is generated to update the controller

→ generate solutions with high weighted accuracy & guaranteed hardware specification

## ■ NASAIC: CONTROLLER AND SELECTOR



## ■ NASAIC: EVALUATOR



### Evaluator: neural architecture accuracy

**Given:**

- An **identified** neural network architecture  $D_i$  for each task
- A **held-out validation dataset** for task  $T_i$  in a total of  $|T|$  tasks

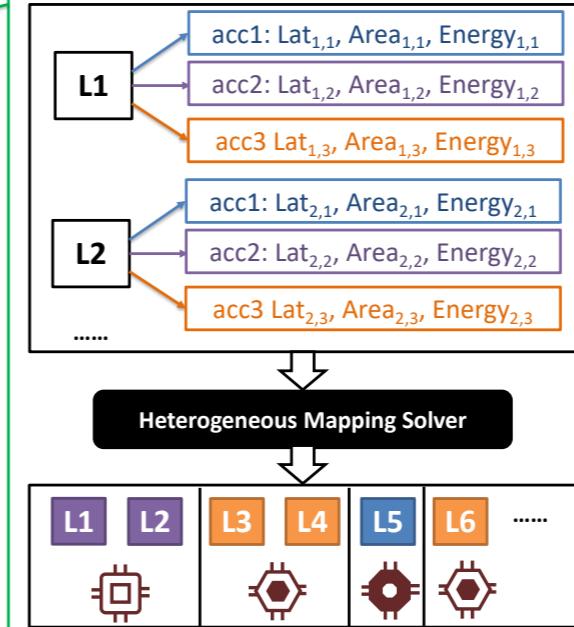
**Do:**

- Training and validation to obtain accuracy  $acc_i$  of  $D_i$
- Feedback **weighted accuracy**, given weight  $\alpha_i$  of  $D_i$

$$weighted(D) = \sum_{i=1,2,\dots,|T|} \{\alpha_i \times acc_i\}$$

### Evaluator: hardware performance

**Cost Model: MAESTRO**



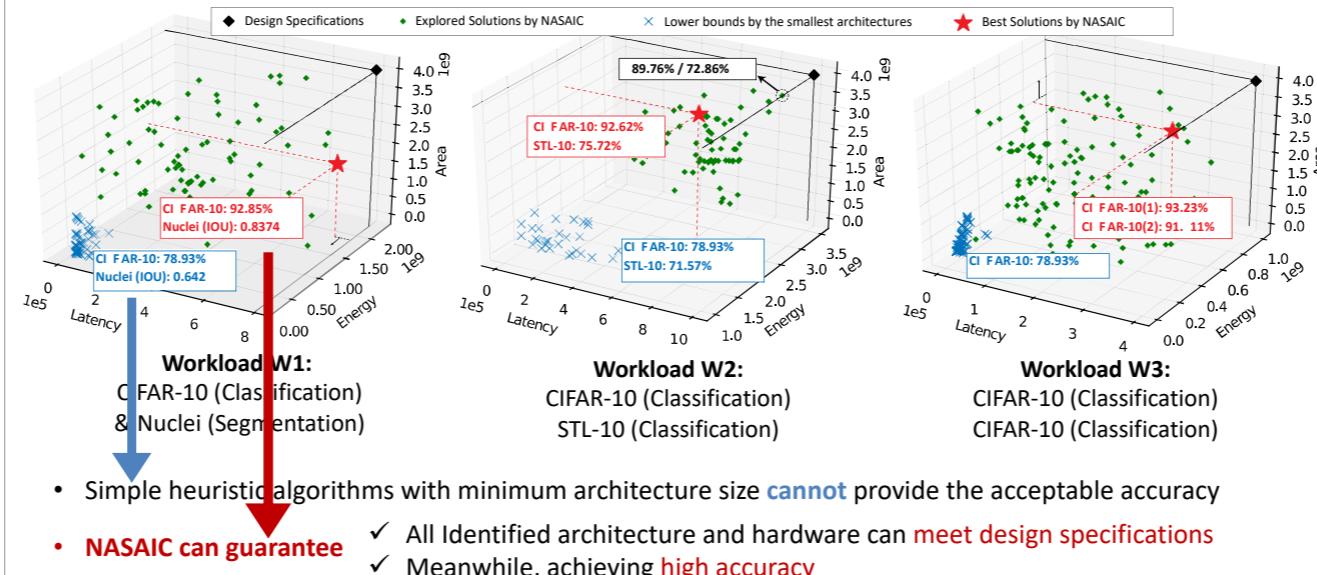
Evaluator uses two passes to train and validate with a hyperparameter based on the accuracy and then calculated in the reward.

For area, need hardware metric. Use state-of-the-art cost model (MAESTRO).

## ■ EXPERIMENT SETUP

- ❖ **Application workloads:** CIFAR-10, STL-10
- ❖ **Backbone architectures:** ResNet9, CIFAR-10, STL-10,
- ❖ **Hardware configuration:** PEs as 4096 and the maximum NoC bandwidth as 64GB/s
- ❖ **ASICNAS setting:** explore the search space for 500 episodes and explore 10 accelerator designs in each episode.

## ■ RESULTS: DESIGN SPACE EXPLORATION



58

Data movement bottleneck with prevent the accelerator from getting higher accuracy

## ■ COMPARISON RESULTS ON MULTI-DATASET WORKLOADS

Workload	Sequential NAS and ASIC Design (Best Accuracy)			CI FAR-10 Nuclei	L/ cycles	E/ nJ	A/ $\mu\text{m}^2$
	NAS → ASIC	dl a, 2112 , 48	1. Dataflow 2. Num PE 3. Bandwidth				
W1 CIFAR-10 Nuclei	ASIC → HW-NAS	dl a, 576 , 56 792 , 8	CI FAR-10 Nuclei	91.98% 83.72%	5.8e5 ○	1.94e9 ○	3.82e9 ○
	NASAIC	dl a, 68 , 56 shi, 1728 , 8	CI FAR-10 STL-10	92.85% 83.74%	7.77e5 ○	1.43e9 ○	2.03e9 ○
	Ours	dl a, 2112 , 24 shi, 1536 , 40	CI FAR-10 STL-10	94.17% 76.50%	9.31e5 ○	3.55e9 ×	4.83e9 ×
W2 CIFAR-10 STL-10	ASIC → HW-NAS	dl a, 2112 , 24 shi, 1184 , 24	CI FAR-10 STL-10	92.53% 72.07%	9.69e5 ○	2.90e9 ○	3.86e9 ○
	NASAIC	dl a, 2112 , 40 shi, 1184 , 24	CI FAR-10 STL-10	92.62% 75.72%	6.48e5 ○	2.50e9 ○	3.34e9 ○

×: Violate the design specs.      ○: Meet the design specs.

NAS --> ASIC:

- Highest accuracy
- Cannot meet design spec.

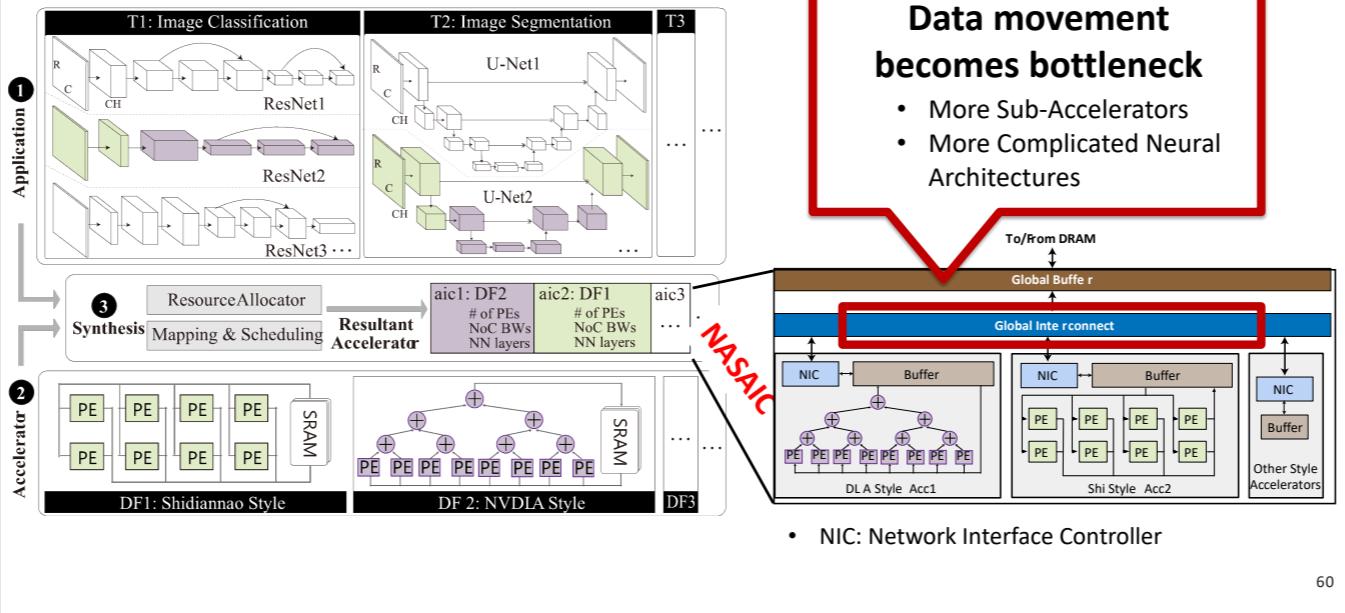
ASIC --> HW-NAS:

- Hardware is not optimized
- Accuracy is low

NASAIC:

- ✓ Best tradeoff between accuracy and hardware efficiency

## ■ PROBLEM AND CHALLENGES

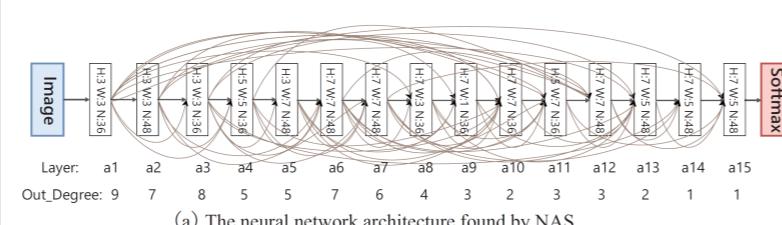


Interconnection between sub-accelerators becomes the challenge when adding to large-scale computation systems.

## OUTLINE

- Hardware Design Space
- Software Design Space
- First Step: Co-Exploration of Neural Architecture and FPGAs
- Second Step: Co-Exploration of Neural Architecture and Heterogenous ASICs
- Third Step: Co-Exploration of Neural Architecture and Scalable ASICs 
- Future Work

## MOTIVATION: SCALABLE NETWORK-ON-CHIP (NOC) FOR NN



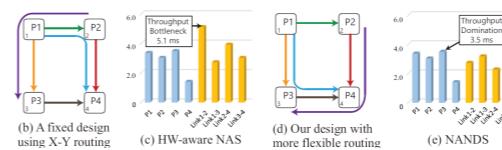
Operation Time	Platforms	Single Processing Element	4 Processing Elements	
			Bus Interconnection	2-D Mesh NoC
Computation (ms)		12.4	3.4	3.4
Data transmission (ms)		—	14.7	6.2

(b) The timing performance of network implementations on different platforms

REF: Zoph, Barret, and Quoc V. Le. "Neural architecture search with reinforcement learning." *ICLR 2017*

### Observations:

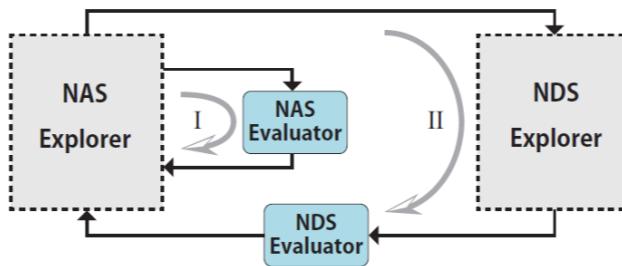
- More PEs: better performance
- Communication: performance bottleneck
- Fixed design: lower performance



65

Increasing number of processing elements creates a communications (interconnection) bottleneck. Network on Chip provides better performance through more flexible routing.

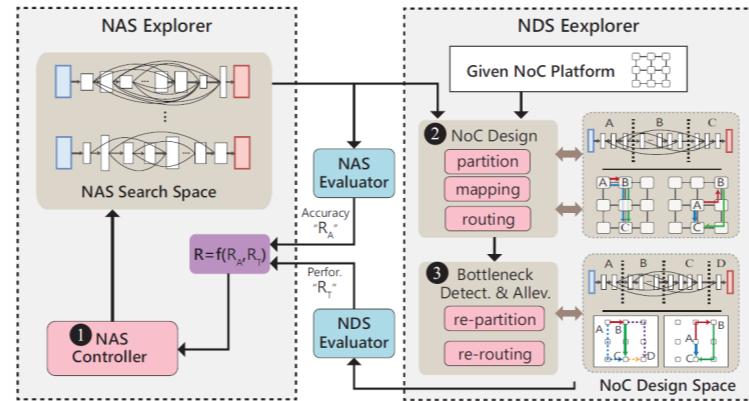
## ■ NANDS: MULTI-PHASE FRAMEWORK



Two exploration loops in NANDS:

- Loop I: Neural Architecture Search.
- Loop II: Automatic Hardware Design.

## ■ NANDS: MULTI-PHASE FRAMEWORK



- ① **NAS Controller:** predict hyperparameters to generate the child neural networks
- ② **NoC Design:** generate hardware design (e.g., partition, mapping and routing) for the input network on a given NoC
- ③ **Bottleneck Detection and Alleviation:** maximize the throughput of NoC.

## ■ NANDS: MULTI-PHASE FRAMEWORK

### ▫ Three kinds of throughput (TP) bottlenecks

- B1: TP is determined by a processing element;
- B2: TP is determined by a NoC link, and the link is occupied by only one data transmission path;
- B3: TP is determined by a NoC link, where there are multiple routing paths will go through

---

#### Algorithm 1 Bottleneck Alleviation

**Input:** (1) NoC with  $PE$  and  $LINK$ , (2) TP, (3) Bottleneck type, (4) latency function  $L$ , (5) partition  $P$ , (6) mapping  $M$ :  $p$  to  $pe$ , (7) routine path  $R$ .

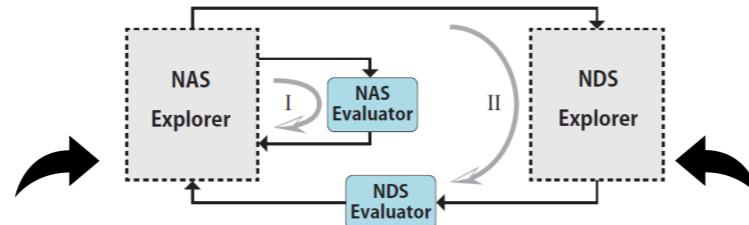
**Output:** A NoC design.

```
1: if Bottleneck is B1, and TP is determined by  $pe_k$ :
2:   Get partition  $p_i$ , s.t.,  $M(p_i) = pe_k$ ;
3:   if  $p_i$  has only 1 layer:
4:     Cannot remove the bottleneck and terminate;
5:   else if NoC has available processing element:
6:     Partition  $p_i \rightarrow p_{i1} + p_{i2}$  to minimize their max latency;
7:   else:
8:     Find  $p_i$ 's neighbor  $p_k$  with the minimum latency;
9:     Move layers in  $p_i$  to  $p_k$  to minimize  $\max(L_{M(p_i)}, L_{M(p_k)})$ ;
10:  else if Bottleneck is B2, and the only routing path is  $pe_i \rightarrow pe_j$ :
11:    Merge partitions on  $pe_i$  and  $pe_j$  to hide communication;
12:  else if Bottleneck is B3, and the link is  $lk$ :
13:    Obtain  $pe_i \rightarrow pe_j$  passing  $lk$  with the maximum data volume;
14:  Re-routing  $pe_i \rightarrow pe_j$  to detour at  $lk$ ;
```

---

## ■ EXPERIMENT SETUP

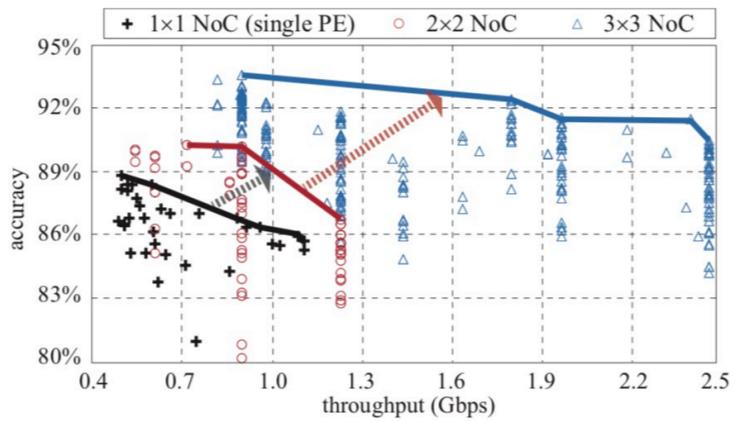
❖ Datasets: CIFAR-10, CIFAR-100, and STL-10



❖ NAS Space: ResNet as the backbone

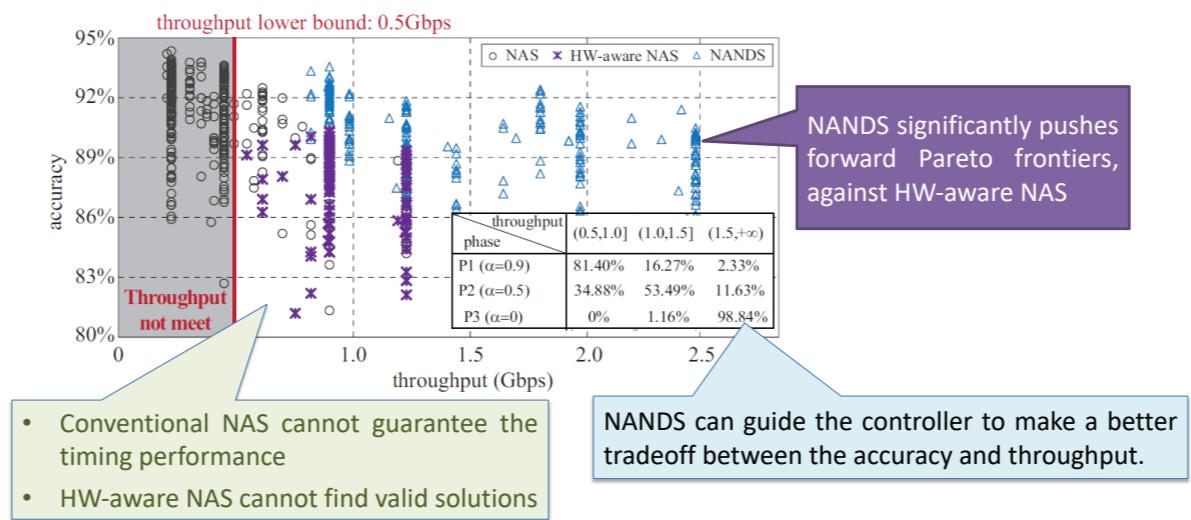
❖ NDS Space:  $2 \times 2$  and  $3 \times 3$  2-D Mesh NoCs

## ■ EXPERIMENTAL RESULTS



Pareto frontiers of accuracy-throughput tradeoffs captured by NANDS on CIFAR-10 can be significantly pushed forward with increasing NoC size.

## ■ EXPERIMENTAL RESULTS



## ■ EXPERIMENTAL RESULTS

TABLE I  
COMPARISON OF THE SEARCH TIME BETWEEN PURE NDS, NAS, HW-AWARE NAS, AND NANDS, ON THREE COMMON DATASETS

Dataset	Spec. (Gbps)	Models	Arch. Property			Accuracy		Throughput			Elapsed Time	
			Depth	Para. ( $\times 10^6$ )	MACs (GOP)	(%)	degr.	(Gbps)	Sat.	impr.	(minute)	impr.
CIFAR-10	0.50	NAS	9	0.89	0.70	94.41%	0.00%	0.22	✗	baseline	1115	baseline
		HW-Aware NAS	8	0.19	0.05	90.95%	-3.46%	0.66	✓	2.94×	164	6.80×
		NANDS (Opt TP)	<b>8</b>	<b>0.20</b>	<b>0.06</b>	<b>91.58%</b>	<b>-2.83%</b>	<b>2.40</b>	✓	<b>10.66×</b>	<b>361</b>	<b>3.09×</b>
		NANDS (Opt Acc.)	<b>10</b>	<b>0.40</b>	<b>0.21</b>	<b>93.59%</b>	<b>-0.82%</b>	<b>0.90</b>	✓	<b>4.00×</b>		
CIFAR-100	0.45	NAS	12	1.04	1.02	76.58%	0.00%	0.22	✗	baseline	1863	baseline
		HW-Aware NAS	8	0.19	0.07	71.43%	-5.15%	0.28*	✗	1.25×	246	7.57×
		NANDS (Opt TP)	<b>8</b>	<b>0.25</b>	<b>0.15</b>	<b>72.22%</b>	<b>-4.36%</b>	<b>0.90</b>	✓	4.00×		
		NANDS (Opt Acc.)	<b>12</b>	<b>0.63</b>	<b>0.46</b>	<b>75.58%</b>	<b>-1.00%</b>	<b>0.45</b>	✓	2.00×	<b>594</b>	<b>3.14×</b>
STL-10	0.6	NAS	11	2.95	2.13	76.45%	0.00%	0.45	✗	baseline	2928	baseline
		HW-Aware NAS	12	1.70	0.50	74.25%	-2.20%	0.61	✓	1.25×	402	7.28×
		NANDS (Opt TP)	<b>11</b>	<b>2.02</b>	<b>1.02</b>	<b>75.83%</b>	<b>-0.62%</b>	1.07	✓	2.37×		
		NANDS (Opt Acc.)	<b>13</b>	<b>2.65</b>	<b>1.45</b>	<b>76.45%</b>	<b>0.00%</b>	0.60	✓	1.32×	<b>1059</b>	<b>2.76×</b>

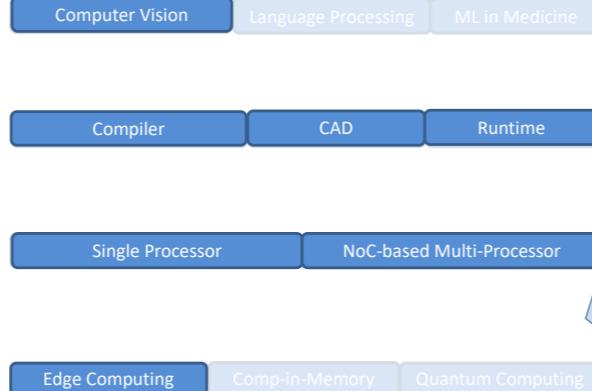
“\*”: relax spec., HW-aware NAS cannot guarantee throughput of 0.45Gbps.

## OUTLINE

- Hardware Design Space
- Software Design Space
- First Step: Co-Exploration of Neural Architecture and FPGAs
- Second Step: Co-Exploration of Neural Architecture and Heterogenous ASICs
- Third Step: Co-Exploration of Neural Architecture and Scalable ASICs
- Future Work 

## ➤ EMERGING HARDWARE & APPLICATION CREATE NEW OPPORTUNITIES!

Software



Hardware

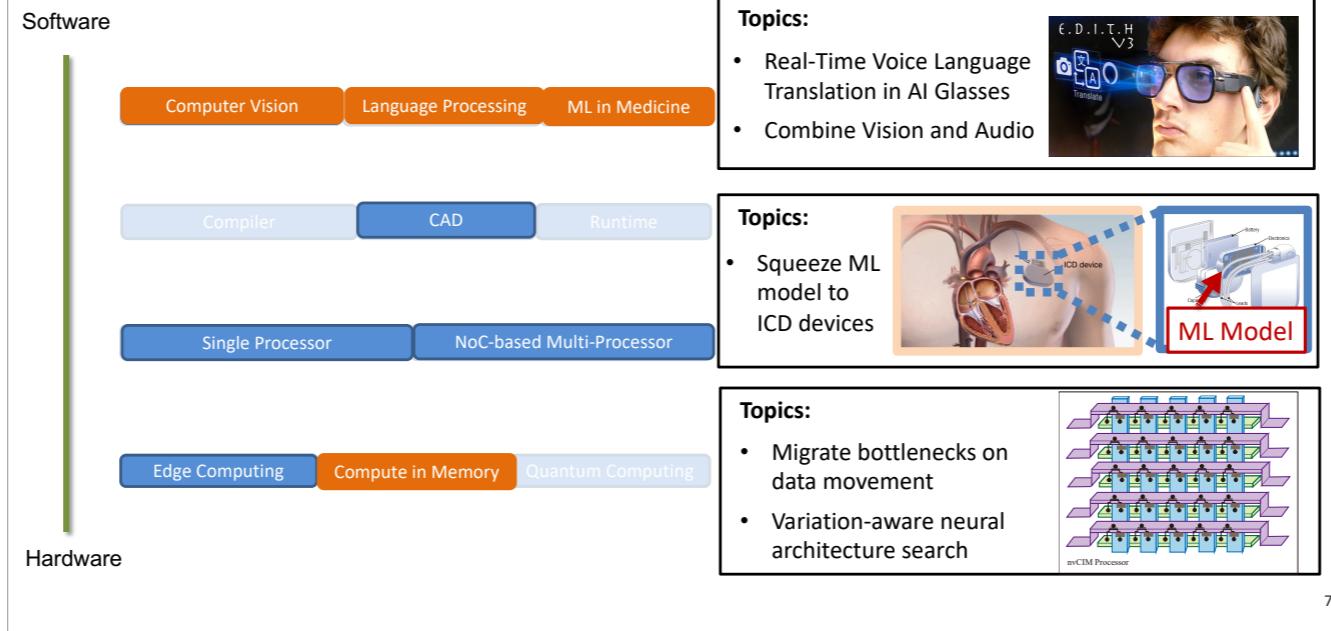
### Co-Explore NAS & Edge Computing (Postdoc work):

- ACM TECS'19
- CODES+ISSS'19 (BPN)
- FPGA'19
- DAC'19 (BPN)
- IEEE TCAD'18
- CASES'18
- DAC'20
- ASP-DAC'19 (BPN)

### Computing Architecture (Ph.D. work):

- ASP-DAC '19 (BPN)
- IEEE TC '18
- IEEE TCAD '18
- CODES+ISSS '18
- FGCS '17
- IEEE TPDS '17
- ASP-DAC '17
- DAC '17
- EMSOFT'17
- IEEE TVLSI '16
- ASP-DAC'16 (BPN)
- HPCC'15
- ISVLSI'14
- RTCSA '14

## ➤ EXPAND THE HW-SW CO-EXPLORATION



75

## ➤ CO-EXPLORE QUANTUM COMPUTING

### Software

Machine Learning Applications

Compiler

CAD

Runtime

Single Processor

NoC-based Multi-Processor

Edge Computing

Comp-in-Memory

Quantum Computing

### Topics:

- Co-Explore Neural Network with Quantum computer



IBM & University of Notre Dame Quantum program

(access IBM Q 53 qbits)



### Hardware